

Reasoning based on symbolic and parametric knowledge bases: a survey

Mayi Xu^{1, †}, Yunfeng Ning^{1, †}, Yongqi Li¹, Jianhao Chen¹,
Jintao Wen¹, Yao Xiao¹, Shen Zhou¹, Birong Pan¹, Zepeng Bao¹,
Xin Miao¹, Hankun Kang¹, Ke Sun¹ & Tiejun Qian^{1, *}

School of Computer Science, Wuhan University, Wuhan 430072, China

Abstract Reasoning is fundamental to human intelligence, and critical for problem-solving, decision-making, and critical thinking. Reasoning refers to drawing new conclusions based on existing knowledge, which can support various applications like clinical diagnosis, basic education, and financial analysis. Though a good number of surveys have been proposed for reviewing reasoning-related methods, none of them has systematically investigated these methods from the viewpoint of their dependent knowledge base. Both the scenarios to which the knowledge bases are applied and their storage formats are significantly different. Hence, investigating reasoning methods from the knowledge base perspective helps us better understand the challenges and future directions. To fill this gap, this paper first classifies the knowledge base into symbolic and parametric ones. The former explicitly stores information in human-readable symbols, and the latter implicitly encodes knowledge within parameters. Then, we provide a comprehensive overview of reasoning methods using symbolic knowledge bases, parametric knowledge bases, and both of them. Finally, we identify the future direction toward enhancing reasoning capabilities to bridge the gap between human and machine intelligence.

Keywords Reasoning, symbolic knowledge base, parametric knowledge base, pre-trained language models, knowledge graphs

Citation Title for citation. Sci China Inf Sci, for review

1 Introduction

Reasoning refers to inferring new conclusions from existing knowledge [1], which is fundamental to human intelligence and essential for complex tasks such as problem-solving, decision-making, and critical thinking. The cognitive process of reasoning involves using evidence, arguments, and logic to draw conclusions or make judgments [2], which can provide back-end support for various real-world applications, such as clinical diagnosis [3–5], basic education [6–8], and financial analysis [9–11]. Reasoning ability is central to human intelligence, yet modern natural language processing systems still struggle to reason based on the information they are given or have already learned [12–15]. The study of reasoning is essential in fields like neurosciences [16], psychology [17], philosophy [18,19], and computer science [20], as it helps to narrow the gap between human and machine intelligence [15]. Hence, building an artificial intelligence system capable of reasoning is both the goal of the research community and the way to improve the performance of complex applications [1].

With the rapid development of reasoning technology, some surveys [1, 2, 15, 21–30] summarized the reasoning methods from different perspectives. For instance, there was a survey [1] researching reasoning using natural language format, including classical logical reasoning, natural language inference, multi-hop question answering, and commonsense reasoning. A few of them [22–25] emphasized the reasoning based on the structured facts in knowledge graphs (KGs), like temporal knowledge graph reasoning and multi-modal knowledge graph reasoning. Some studies [26–30] paid attention to the knowledge sources that the reasoning methods used for question answering, e.g., Wikidata KG and Wikipedia corpus. More recent surveys [2, 15, 21] summarized the reasoning methods by prompting large language models, such as chain-of-thought series and self-reflection series methods.

[†]Co-first author

* Corresponding author (email: qty@whu.edu.cn)

Despite the valuable perspectives provided by these surveys, none of them has summarized reasoning methods from the view of their dependent knowledge base. As previous work [1, 2, 31] points out, reasoning is a process of integrating multiple existing knowledge to derive some new conclusions about the world, and current reasoning methods rely heavily on knowledge bases. However, both the scenarios to which the knowledge bases are applied and their storage formats are significantly different. Hence, investigating reasoning from the perspective of the knowledge base helps us gain a deeper understanding of the challenges and future directions.

In this paper, we review related work by focusing on the underlying knowledge base that supports reasoning methods. We begin by classifying these knowledge bases into two types based on their storage formats: symbolic and parametric, where symbolic knowledge bases present information in human-readable symbols like KGs and tables, and parametric ones encode information implicitly within parameters. Then, we investigate the reasoning methods based on symbolic knowledge bases, parametric knowledge bases, and both of them, respectively. Finally, we explore the challenges and potential future directions for reasoning with both symbolic and parametric knowledge. The overall framework is shown in Figure 1.

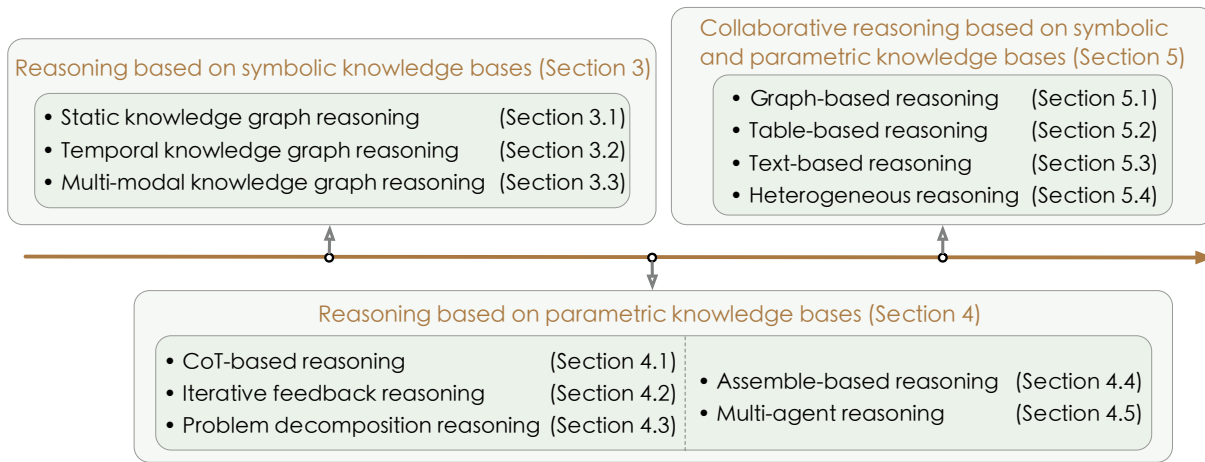


Figure 1 Overall framework of reasoning-related methods.

In summary, the main contributions of this survey are:

- To the best of our knowledge, we are the first to provide a comprehensive survey on reasoning studies from the perspective of their dependent knowledge bases.
- We conduct a thorough investigation of various types of reasoning methods that utilize symbolic knowledge bases, parametric knowledge bases, and their combination, whereas previous reviews only focused on one of them.
- We have meticulously summarized the challenges and future research directions related to reasoning, which will contribute to advancing the development of this field.

Organization of this survey: we first introduce the background in Section 2. Then, we systematically introduce different reasoning tasks in Section 3, 4, 5. We discuss the challenges and future research directions in Section 6. Finally, we conclude this paper in Section 7.

2 Background

In this section, we first introduce the concepts of symbolic and parametric knowledge bases. In artificial intelligence, the symbolic knowledge bases and parametric knowledge bases represent different paradigms of knowledge representations that align with symbolism [32] and connectionism [33], respectively. The symbolic knowledge base involves explicit knowledge and logical structures for reasoning. It is also fundamental to symbolic AI, which focuses on rule-based manipulation of symbols [34]. In contrast, the parametric knowledge base is associated with connectionism, where neural networks capture knowledge implicitly through learned parameters, emphasizing adaptability and pattern recognition [35]. Finally, we will introduce the taxonomy of reasoning in detail.

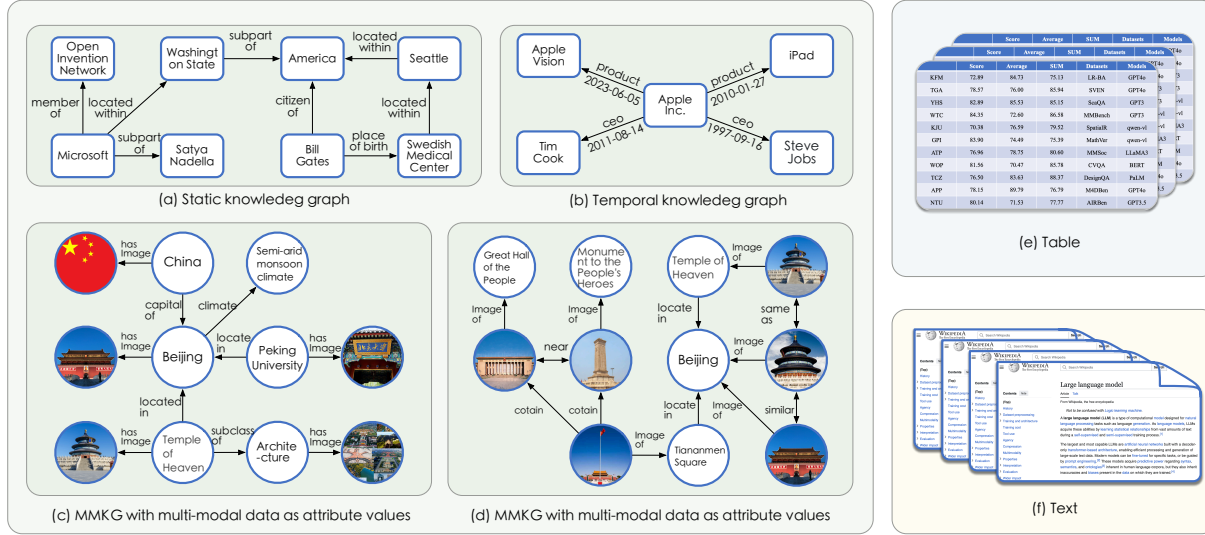


Figure 2 Different symbolic knowledge bases.

2.1 Symbolic knowledge bases

The symbolic knowledge bases contain KGs, tables, and text, the first two are structured and the last one is unstructured. The structured KGs can be partitioned into static knowledge graphs (SKGs), temporal knowledge graphs (TKGs), and multi-modal knowledge graphs (MMKGs). The following part of this subsection introduces the details of these symbolic knowledge bases.

Static knowledge graph: As shown in Figure 2 (a), an SKG is a structured semantic knowledge base that can express various associations between entities in a graphical manner [36]. An SKG contains many factual triplets (h, r, t) , where h , r , and t represent the head entity, the relation, and the tail entity, respectively. For example, $(Microsoft, located\ within, Washington\ State)$ represents “Microsoft is located within Washington State”. An SKG is represented as $(\mathcal{E}, \mathcal{R}, \mathcal{F})$, where \mathcal{E} , \mathcal{R} , and \mathcal{F} denote the entity set, the relation set, and the triplet set $\{(h, r, t)\} \subseteq \mathcal{E} \times \mathcal{R} \times \mathcal{E}$, respectively.

Temporal knowledge graph: As shown in Figure 2 (b), TKGs can store temporal information such as events that evolve over time. The quadruples $(h, r, t, time)$ are the basic unit of TKGs, where h , r , t , and $time$ represent the head entity, the relation, the tail entity, and the timestamp, respectively. For example, the quadruple $(Apple\ Inc., product, iPad, 2010-01-27)$ means “On January 27, 2012, Apple Inc. produced iPad”. A TKG is represented as $(\mathcal{E}, \mathcal{R}, \mathcal{T}, \mathcal{F})$, where \mathcal{E} , \mathcal{R} , \mathcal{T} , \mathcal{F} denote the entity set, relation set, the timestamp set, and the quadruple set $\{(h, r, t, time)\} \subseteq \mathcal{E} \times \mathcal{R} \times \mathcal{E} \times \mathcal{T}$, respectively. Another representation way of TKGs is $(G_{t_1}, G_{t_2}, \dots, G_{t_n})$, where G_{t_i} denotes the SKG containing all triplets happened in timestamp t_i .

Multi-modal knowledge graph: The MMKG integrates various multi-modal data into one SKG, such as text, images, and audio. As shown in Figure 2 (c) and (d), MMKGs are generally divided into two types [36]. One represents multi-modal data as new entities, and the other one describes them as entities’ attributes.

Structured table: As shown in Figure 2 (e), a structured table contains n records and m attributes. Each record can be represented as a vector $\mathbf{r}_i = (a_{i1}, a_{i2}, \dots, a_{im})$, where a_{ij} denotes the value of the j -th attribute in the i -th record. The entire table can be viewed as a set of n records $\{\mathbf{r}_1, \mathbf{r}_2, \dots, \mathbf{r}_n\}$. The format of table makes data easy to store, manage, and analyze [37]. Hence, it is widely used in various fields such as financial statements, scientific research data, and inventory management [38, 39].

Unstructured text: As shown in Figure 2 (f), unstructured text primarily refers the text information without a standardized format or organization, such as books and articles. Unstructured text can also store knowledge and be widely used in various applications with the help of current large language models.

2.2 Parametric knowledge bases

Parametric knowledge bases mainly refer to pre-trained language models (PLMs). PLMs are pre-trained on large-scale text corpora via self-supervised learning and store generalized knowledge in pa-

rameters, enabling them to reason with implicit knowledge in parameters [40]. Based on architecture differences, current PLMs can be divided into encoder-only, decoder-only, and encoder-decoder.

Encoder-only PLMs: The representative encoder-only PLMs mainly include BERT [41] and its variants, like RoBERTa [42], ALBERT [43], DeBERTa [44], XLM [45], XLNet [46], UNILM [47], and ELECTRA [48]. Different pre-training strategies are designed to incorporate pre-training knowledge into their parameters. For instance, BERT is pre-trained through the application of mask language modeling. RoBERTa, which is built upon and optimized from BERT, makes use of dynamic masking. In terms of ELECTRA, it introduces the replaced token detection pre-training task, which enhances the learning efficiency and performance of the model. To effectively utilize parametric knowledge in encoder-only PLMs, a direct approach is extracting semantic representations from the input text. This method is widely used in reasoning tasks, such as open-domain question answering [49–52] and machine reading comprehension [53–56], to incorporate pre-trained knowledge.

Decoder-only PLMs: Pre-trained language models adopting the decoder-only architecture primarily focus on text generation tasks. Based on the autoregressive attribute, these models generate coherent texts by successively predicting the forthcoming content based on the already generated text. Representative decoder-only PLMs encompass GPT [57], LLaMA [58], Qwen [59], and Mistral [60]. Recently, decoder-only PLMs have demonstrated remarkable potential within the domain of text generation. Consequently, most PLMs with large parameter scales (also known as large language models, LLMs) use decoder-only architecture. During reasoning, we can elicit knowledge from LLMs through prompting methods such as chain of thought [61], iterative feedback [62], problem decomposition [63], assemble [64] and multi-agent [65]. In contrast, when using small language models (SLMs) for reasoning, it may necessary to fine-tune them using task-specific or domain-specific data to leverage the parametric knowledge contained within them effectively [66,67].

Encoder-decoder PLMs: Encoder-decoder PLMs integrate both encoder and decoder components. It uses the encoder to model input text features and the decoder for text generation, with the advantage of separating vector spaces for text understanding and generation. However, it has some limitations such as high complexity and significant demands on pre-training time and computational resources. Representative encoder-decoder PLMs include T5 [68], BART [69], mBART [70], MASS [71], and ChatGLM [72]. Thanks to its structural advantages, we can utilize the encoder’s semantic vector features and the decoder’s target text for knowledge reasoning [73].

2.3 Taxonomy of reasoning

From the perspective of the dependent knowledge bases, reasoning methods can be divided into three categories: reasoning methods based on symbolic knowledge bases, reasoning methods based on parametric knowledge bases, and collaborative reasoning methods based on symbolic and parametric knowledge bases.

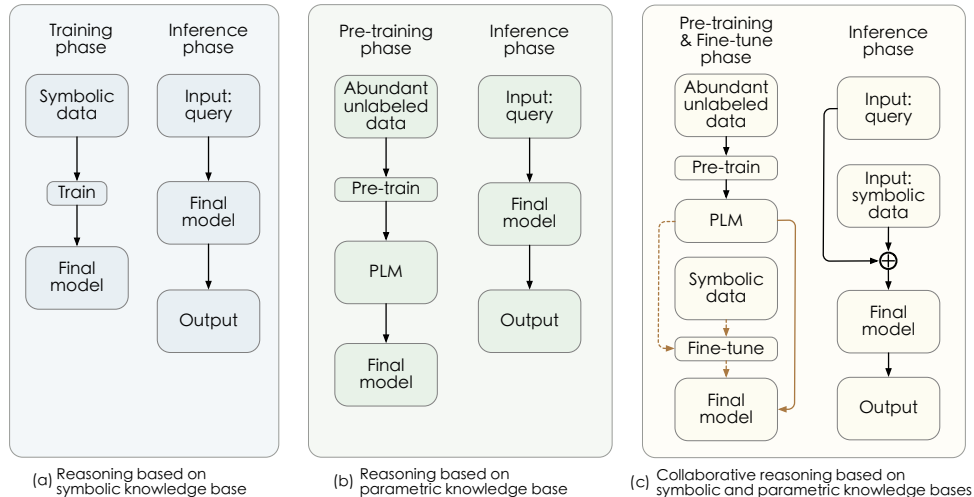


Figure 3 Different reasoning methods.

As shown in Figure 3 (a), reasoning methods based on symbolic knowledge bases train a model to learn the knowledge in the symbolic knowledge bases during the training phase. During inference, only a query is inputted to the model. The model’s reasoning capability comes from the knowledge modeled in the training phase. Typical tasks for this type of method, such as static knowledge graph reasoning, require reasoning for an unknown element in an incomplete triplet or quadruple. [74, 75]

As shown in Figure 3 (b), reasoning methods based on parametric knowledge bases leverage the knowledge in PLMs’ parameters. During inference, only a query is inputted. The model’s reasoning capability comes from the parametric knowledge modeled during pre-training. Compared to the training phase in Figure 3 (a), the pre-training phase is usually task-independent. Typical tasks for this type of method, such as mathematical reasoning, require reasoning for a natural language answer given a knowledge-intensive question [61–63, 76, 77].

As shown in Figure 3 (c), collaborative reasoning methods based on symbolic and parametric knowledge bases also model parametric knowledge during the pre-training phase. Some methods further fine-tune the PLMs using knowledge from symbolic knowledge bases to enhance the domain knowledge learning. During inference, these methods retrieve relevant knowledge from the symbolic knowledge bases and combine it with the parametric knowledge in PLMs to enhance reasoning performance. Similar to reasoning methods based on parametric knowledge bases, this type of method, such as knowledge graph questions answering, also performs reasoning in the form of natural language question answering [78–81].

3 Reasoning based on symbolic knowledge base

In this section, we investigate the reasoning methods based on symbolic knowledge bases. Based on the structural types of KGs, we investigate static knowledge graph reasoning, temporal knowledge graph reasoning, and multi-modal knowledge graph reasoning. The overall taxonomy of reasoning methods based on symbolic knowledge bases is shown in Figure 4.

3.1 Static knowledge graph reasoning

Static knowledge graph reasoning refers to the completion of incomplete triplets (incomplete knowledge) based on the given fact triplets (existing knowledge) in the SKG, thereby obtaining new complete factual triplets (new knowledge). According to the differences in query form and output form, static knowledge graph reasoning tasks can be divided into three sub-tasks: traditional SKG reasoning, multi-hop SKG reasoning, and SKG complex logical query answering. In general, multi-hop SKG reasoning requires providing an explicit reasoning path while completing the triplets. SKG complex logical query answering requires modeling the complex logical symbols.

3.1.1 Traditional SKG reasoning

Task definition: Given the query $(h, r, ?)$, $(?, r, t)$, or $(h, ?, t)$, traditional SKG reasoning methods aim to predict the missing entity or relation directly. Generally, the traditional SKG reasoning methods first learn the representations of relations and entities in SKGs. Then, it constructs a scoring function to calculate the validity of possible triples. According to the representation approaches, traditional SKG reasoning methods can be categorized into three types: translation-based, tensor decompositional, and neural-based [36, 163].

3.1.1.1 Translation-based methods

Translation-based methods map the entities and relations to a vector space, where entities with similar semantics are close in the distance. The role of the relation is to project the head entity representation onto the tail entity representation. A common approach is to represent relations as translation vectors from the head entity to the tail entity.

The first translation-based method TransE [74] regards the relation as a translation operation, that is, $h + r \approx t$, where h , r , and t belong to the same vector feature space. However, it cannot handle certain specific relations, such as one-to-many, many-to-one, symmetric, and transitive relations. To adapt to the multiple meanings of entities and relations, TransH [82] interprets relations as transformation operations on hyperplanes, TransR [83] models different atoms in the SKG in different vector spaces with the projection matrix M_r for each relation r , and TransD [84] dynamically constructs mapping matrices for

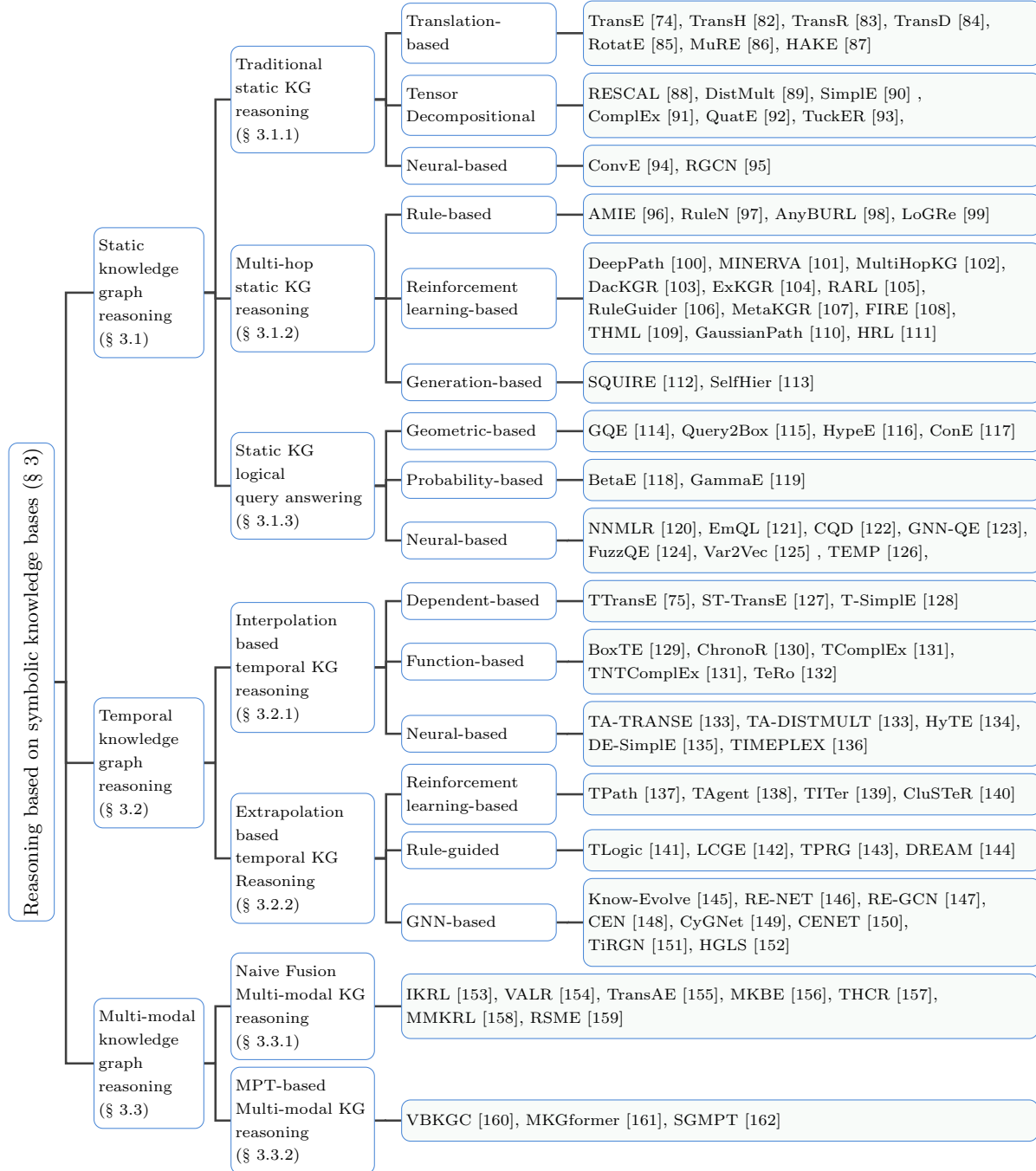


Figure 4 Taxonomy of reasoning methods based on symbolic knowledge base.

every entity-relation pair. Furthermore, some methods propose vector spaces that differ from traditional approaches. For example, RotatE [85] considers relation as rotation in complex space, MuRE [86] utilizes the Poincaré ball model in hyperbolic space, and HAKE [87] models semantic hierarchies based on polar coordinate space rather than relational patterns.

3.1.1.2 Tensor decompositional methods

Tensor decompositional methods focus on capturing the pairwise interaction between atoms in an SKG to exploit the similarity of the latent features. For example, RESCAL [88] models entities as vectors and relations as matrices to capture the interaction amounts between the elements at corresponding positions of the head and tail entity vectors. However, RESCAL is high in computational complexity and cannot model asymmetric relations. Therefore, many methods like DistMult [89], Simple [90], ComplEx [91], QuatE [92], TuckER [93] utilize tensor decomposition to better model the semantic similarity of triples and reduce the computational complexity of the model.

3.1.1.3 Neural-based methods

Translation-based and tensor decompositional methods create a vector or matrix for every entity and relation, which requires large storage space for large-scale SKGs. To address the challenge, neural network methods utilize multi-layer neural networks to learn the representations of entities and relations, which avoids the storage of vectors or matrices. For example, ConvE [94] presents a simple multi-layer convolutional architecture. Inputting the head entity and relation, the encoder of ConvE calculates their deep features, and the decoder of ConvE maps them into the entity space to predict the target tail entity. Besides, since graph neural network (GNN) [164] has proven to be effective in mining the structural information of graph, it is widely applied for traditional SKG reasoning. For example, RGCN [95] encodes entities by aggregating neighbor information instead of optimizing with a single entity vector, resulting in significant performance improvements.

3.1.2 Multi-hop SKG reasoning

Task definition: Given a query $(h, r, ?)$, multi-hop SKG reasoning methods aim to predict the target tail entity t through a n -hop reasoning path $\tau : h \xrightarrow{r_1} e_1 \xrightarrow{r_2} e_2 \cdots \xrightarrow{r_n} e_n$, where e_i and r_i represent the entity and the relation in the path τ . The last entity e_n in τ is treated as the predicted target tail entity t . Compared to traditional SKG reasoning methods, multi-hop SKG reasoning methods provide better interpretability. The multi-hop SKG reasoning methods can be categorized into rule-based, reinforcement learning-based (RL-based), and generation-based.

3.1.2.1 Rule-based methods

Rule-based methods automatically induce logical rules from SKG and predict missing entities or relations by matching queries to the rules, typically without model training. The rules exist in the form of symbolic chains, such as $mother_of(m, c) \wedge married_to(m, f) \Rightarrow father_of(f, c)$, in which some facts can be inferred from other facts. Some rule-based methods such as AMIE [96], RuleN [97], AnyBURL [98], LoGre [99] aim to efficiently mine and utilize rules in KGs. Although these methods achieve remarkable performance, they are hard to generalize in practice due to the limitation of symbolic representation.

3.1.2.2 Reinforcement learning-based methods

RL-based methods model the multi-hop reasoning process as a markov decision process (MDP), where the agent walks on SKG to search the target entities. DeepPath [100] first adopts the RL framework to search the reasoning paths and target relations given to the head and tail entities. MINERVA [101] proposes a more complex and practical scenario to find the target tail entities given the relations and head entities. Following MINERVA, most RL-based methods are devoted to tackling the sparse rewards problem and trying to design a more efficient policy network in incomplete SKG. For instance, Multi-HopKG [102] is one of the first to use embedding models to estimate the rewards of unobserved targets, thereby reducing the impact of incorrect negative samples. Similar methods include DackGR [103], ExKGR [104], RARL [105], and RuleGuider [106].

Moreover, several RL-based methods have integrated meta-learning to address the decreased reasoning capability of reinforcement learning in few-shot relation scenarios. For instance, Meta-KGR [107] is the

first to apply meta-learning to multi-hop KG reasoning, using the meta-learning algorithm MAML [165] to learn meta-parameters from high-frequency relations and then quickly adapts to sparse relations. FIRE [108] uses heterogeneous neighbor information to enhance entity embeddings and leverages knowledge graph embeddings to compute structural relevance, thereby reducing the search space. THML [109] proposes a difficulty-aware meta-reinforcement learning method that trains difficulty-aware batches to predict missing elements, as well as a two-level difficulty-aware sampling strategy to effectively generate new difficulty-aware batches, greatly enhancing the generalization capability.

Recently, there have been some quite creative methods. For example, GaussianPath [110] points out that agents in traditional RL-based methods are prone to get trapped in reasoning paths. It proposes a bayesian multi-hop reasoning paradigm to capture the uncertainty of reasoning paths to explore a broader range of reasoning paths. Besides, inspired by how humans handle ambiguous situations, HRL [111] proposes a high-level policy to learn historical information and a low-level policy to recognize relation clusters. It decomposes the complex action space to express the multiple semantics of relations.

3.1.2.3 *Generation-based methods*

Generation-based methods adopt a generative framework to generate the reasoning paths step by step. For instance, SQUIRE [112] utilizes an encoder-decoder model to translate the query to a reasoning path. By leveraging the rule-enhanced and iterative training strategy, SQUIRE performs better than rule-based and RL-based methods. In cold-start multi-hop reasoning, the model always lacks precise guidance and explicit paths. To overcome these challenges, SelfHier [113] designs an effective generation-based model to explore the reasoning paths by hierarchical guidance and self-verification strategies. Overall, the generation-based methods not only obtain high performance but also has low time complexity.

3.1.3 *SKG complex logical query answering*

Task definition: Given an incomplete first-order logical (FOL) query, SKG complex logical query answering methods aim to predict an target entity set. For instance, the FOL query $C_? . \exists P : assoc(d_1, P) \wedge assoc(d_2, P) \wedge target(P, C_?)$ means that “identify potential drugs $C_?$ that can act on proteins P associated with the disease d_1 and d_2 ”. SKG complex logical query answering methods first transform the given FOL query into a region in embedding space. Then, they output the entities within the region as the target entity set. Based on the difference of embedding space, the SKG complex logical query answering methods can be categorized into geometric-based, probability-based, and neural-based.

3.1.3.1 *Geometric-based methods*

Geometric-based methods transform queries into geometric regions with clear boundaries. The composition and transformation of queries align with the composition and transformation of their geometric regions. For example, GQE [114] proposes the geometric projection operator and the geometric intersection operator, which embeds basic query into a single point and combines basic queries into complex FOL query, respectively. Query2Box [115] effectively models query graph embeddings through hypergeometric box embeddings, where the query is represented as a box using the center and the boundary offset. All entities that fall within the box region of the query are considered the target entity set. In addition, HypeE [116] learns representations of entities and relations as hyperboloids in a Poincaré ball. ConE [117] represents entities and queries as Cartesian products of two-dimensional cones.

3.1.3.2 *Probability-based methods*

Geometric-based methods necessitate that target entities strictly fall within the region defined by the query. However, this requirement often fails to reflect real-world scenarios where entities exhibit semantic diversity. Probability-based methods address this issue by modeling queries and entities as more flexible probability distributions. For example, BetaE [118] uses probabilistic distributions with bounded support, specifically the Beta distribution, and embeds queries/entities as distributions, which allows it to faithfully model uncertainty. Similarly, GammaE [119] utilizes Gamma distribution to capture more features of entities and queries. Based on the linear property and strong boundary support of Gamma distribution, GammaE effectively avoids generating ambiguous answers.

3.1.3.3 Neural-based methods

Neural-based methods model the correspondence between queries and entities using neural network. For example, NNMLR [120] uses multilayer perceptron (MLP) and MLP-Mixer to model the basic queries and their combinations. It computes the distance between the query and the entities to rank the answers. Similar methods include EmQL [121], CQD [122], GNN-QE [123], FuzzQE [124], Var2Vec [125], and TEMP [126].

3.2 Temporal knowledge graph reasoning

The objective of temporal knowledge graph reasoning is to leverage existing events and knowledge to reason about unseen events or predict future events. Previous research [36] primarily categorizes reasoning tasks into two scenarios: interpolation and extrapolation, depending on whether the model has seen the timestamps in the query. From a reasoning perspective, interpolation generally completes missing facts by analyzing known knowledge in TKGs, while extrapolation focuses on predicting unknown events by learning embeddings of entities and relations from historical facts on continuous TKGs.

3.2.1 Interpolation-based TKG reasoning

Task definition: Given a TKG with facts from $time_0$ to $time_T$, the Interpolation-based TKG reasoning method aims to complete missing quadruple $(h, r, ?, time_i)$ or $(?, r, t, time_i)$ in history ($time_0 \leq time_i \leq time_T$). Interpolation-based TKG reasoning methods can be divided into dependent-based, function-based, and neural-based.

3.2.1.1 Dependent-based methods

Dependent-based methods generally do not involve direct manipulation of timestamps. Instead, they associate each timestamp with the relevant entity or relation, capturing the evolution of entities or relations over time. For example, TTransE [75] extends the traditional TransE [74] model by jointly encoding relations and timestamps within a unified space. Building on TTransE, ST-TransE [127] introduces a specialized time embedding method that constrains the representation learning of entities and relations. However, TTransE and ST-TransE struggle to manage evolving facts over time effectively. In contrast, T-Simple [128] leverages a fourth-order tensor to model interactions within quadruples, improving its ability to capture temporal associations.

3.2.1.2 Function-based methods

Function-based methods use specialized functions to learn embeddings for timestamps or model the temporal evolution of entities and relations. Specifically, BoxTE [129] extends the static BoxE [166] model by incorporating temporal information through a relation-specific transfer matrix, facilitating the exploration of more complex inference patterns over time. ChronoR [130] associates timestamps with relations, considering each relation-timestamp pair as a rotation that maps the head entity to the tail entity. TComplex and TNTComplex [131] extend the third-order tensor to a fourth-order tensor in complex space for reasoning. Notably, TNTComplex assumes that certain facts remain static over time, separating the TKG into temporal and non-temporal components. Similarly, TeRo [132] incorporates timestamps into the embeddings of head and tail entities in complex space to capture their temporal evolution, and represents the relation as a rotation that maps the head entity to the tail entity.

3.2.1.3 Neural-based methods

Neural-based methods typically use convolutional neural networks (CNNs) or long short-term memory (LSTM) networks to encode timestamps. These encoded timestamps help model the evolution of entities and relations by capturing their intrinsic correlations and temporal dependencies. For instance, TA-TRANSE [133] is a temporal-aware version of TransE [74]. It utilizes LSTM to learn time-aware representations of relation, and represents quadruples as a set of triples in the form of (h, r_{seq}, t) , where r_{seq} means relation that may include temporal information with a temporal suffix. Similarly, TA-DISTMULT [133] is a temporal extension of DistMult [89], considering the relation with temporal information as a sequence. Additionally, HyTE [134] is an extension of TransH [82], DE-Simple [135] is an extension of Simple [90]. These methods often consider temporal constraints to enhance temporal reasoning capabilities. For ex-

ample, TIMEPLEX [136] leverages the recurrent nature of certain facts and the temporal interactions between pairs of relations during expansion. These additional temporal constraints can help assess a quadruple’s validity better.

3.2.2 Extrapolation-based TKG reasoning

Task definition: Given a TKG with facts from $time_0$ to $time_T$, the Extrapolation-based TKG reasoning method aims to predict unknown facts $(h, r, ?, time_j)$ or $(?, r, t, time_j)$ that occur in the future ($time_j > time_T$). Extrapolation-based TKG reasoning methods can be categorized into RL-based, rule-guided, and GNN-based.

3.2.2.1 Reinforcement learning-based methods

RL-based methods model path reasoning process through a reinforcement learning framework. For instance, TPath [137] adds the time information as a separate vector to participate in the iteration of environment and agent. TAgent [138] filters the embeddings of candidate actions through a novel gate mechanism based on temporal information to capture temporal evolutionary patterns. TITer [139] defines a relative time encoding function to capture the information of timestamps and designs a time-shaped reward based on the Dirichlet distribution to guide the model’s learning. CluSTeR [140] proposes a beam search strategy to elicit multiple clues from historical facts and uses graph convolutional networks to deduce answers from the clues.

3.2.2.2 Rule-guided methods

Rule-guided methods derive some temporal logic rules from TKGs and utilize them to predict future facts. The rules provide a structured framework to infer logical conclusions, thereby generating more accurate predictions about future states or events. For example, TLogic [141] automatically mines cyclic temporal logical rules by extracting temporal random walks from the graph. LCGE [142] mines the temporal rules with several time constraint patterns to construct a rule-guided predicate embedding regularization strategy for learning the causality among events. Rule-guided methods can also be integrated with RL-based methods to reduce semantic noise during reasoning and enhance the stability of the model. For example, TPRG [143] proposes a similar concept of temporal rules and has made improvements based on TPath [137], achieving certain improvements. DREAM [144] proposes a reinforcement learning framework where the agent can receive adaptive rewards by imitating demonstrations at both the semantic and rule levels to eliminate the issue of sparse rewards.

3.2.2.3 GNN-based methods

Recent advancements in temporal knowledge graphs have leveraged GNNs to manage structural and temporal dependencies. Know-Evolve [145] is a classic and the first temporal knowledge graph reasoning model that models the occurrence of facts (edges) as a multivariate point process over time, thereby learning non-linearly evolving entity representations with a deep recurrent network. As with SKGs, GCN [167] is beneficial for TKGR. For example, RE-NET [146] applies GCN [167] to interpret event occurrences as sequences of subgraphs within TKGs, employing an autoregressive model with a neighborhood aggregation function to enhance interpretability. Both RE-GCN [147] and its advanced version, the CEN [148], holistically treat the entire KG sequence to capture the evolutionary dynamics of entities and relations, effectively pinpointing local historical dependencies. Meanwhile, CyGNet [149] focuses on identifying high-frequency entities by exploiting repetitive patterns in historical data using a copy-generation network. On the other hand, CENET [150] differentiates between historical and non-historical dependencies to better identify entities suited for specific queries. While these models excel at capturing specific long-range facts, they often lack high-order connectivity information and dynamic sequential patterns required for a deeper understanding. To address these shortcomings, TiRGN [151] develop different structural encoders to capture sequential and recurring patterns within historical data. Furthermore, HGLS [152] designs a hierarchical graph framework to model long-term dependencies of entities across different timestamps.

3.3 Multi-modal knowledge graph reasoning

Task definition: The goal of MMKG reasoning methods is similar to SKG reasoning methods, which aims to complete the triplet (h, r, t) when one of h , r , or t is missing. In particular, the entity (h and r) could be text or images or has attributes of text and images. The MMKG reasoning methods can be divided into naive fusion MMKG reasoning and multi-modal pre-trained transformer-based (MPT-based) MMKG reasoning methods.

3.3.1 Naive fusion MMKG reasoning methods

Naive fusion MMKG reasoning methods evolve from traditional SKG reasoning methods. They concentrate on the efficient encoding and integration of multi-modal data. The representation of multi-modal data is achieved by either combining each individual modality’s representation within its own feature space or by projecting different modal representations into a shared latent space.

The earliest work IKRL [153] uses a neural image encoder to construct representations for all images of an entity. Then, the multiple image representations of an entity are combined with the original structure-based representations and trained like TransE [74], thereby learning multi-modal knowledge representations for reasoning. Inspired by IKRL [153], a translation-based Visual and Linguistic Representation Model (VALR) [154] has been proposed, which defines the energy of the triple as the sum of sub-energy functions that leverage both visual, linguistic and structural representations. Similarly, TransAE [155] combines multi-modal autoencoder with TransE [74] model, where the hidden layer of the autoencoder is used to encode multi-modal data. MKBE [156] designs specialized encoding layers, scoring modules, and decoding layers for data in different modals. THCR [157] complements the relational knowledge by learning a shared latent representation that integrates information across those modalities. MMKRL [158] designs a joint learning framework that can be easily extended to any modality and uses an adversarial strategy to enhance its robustness. RSME [159] automatically encourages or filters the influence of visual context to avoid encoding too much irrelevant information.

3.3.2 MPT-based MMKG reasoning methods

MPT-based MMKG reasoning methods use MPT to encode textual and visual features in a unified architecture. Then, it employs graph encoders to integrate structural knowledge with multi-modal knowledge. For example, VBKGC [160] focuses on the co-design of the structural KG model and negative sampling. It consists of an encoding module with VisualBERT [168], a projection module, and a scoring module like TransE [74]. MKGformer [161] utilizes a hybrid transformer architecture with unified input-output and reduces noise from irrelevant images/objects through token-level modal fusion. SGMPT [162] designs a structure-guided fusion module that uses weighted summation and alignment constraint to inject the structural information into both the textual and visual features.

3.4 Datasets

In this section, we have compiled statistics on some commonly-used datasets related to reasoning based on symbolic knowledge bases, including (1) # Ent.: Entity number; (2) # Rel.: Relation number; (3) # Time.: Timestamp number; (4) # Facts: Fact number; (5) Type: Knowledge graph type. In particular, the MMKG is represented by specific combinations of modalities. Taking “KG+TXT+IMG” for example, “KG” means the entity has a simple name or ID, “TXT” means the entity has a textual description as attributes, “IMG” means the entity has single or multiple corresponding images as attributes; (6) Domain: The domain of knowledge stored in the KGs; (7) Source: The source of the KGs, and (8) Links: The storage address of the KGs. The statistical results are shown in Table 1.

4 Reasoning based on parametric knowledge bases

Since reasoning methods based on parametric knowledge bases are often task-independent, this section reviews them from the perspective of their equipped techniques rather than tasks. These methods mainly perform reasoning in the form of question answering.

Datasets	# Ent.	# Rel.	# Time.	# Facts	Type	Domain	Source	Links
ATOMIC [169]	304,388	9	-	785,937	Static	Commonsense	Crowdsourcing	https://
CoDEX_S [170]	2,034	42	-	36,543	Static	General	Wikidata	https://
CoDEX_M [170]	17,050	51	-	206,205	Static	General	Wikidata	https://
CoDEX_L [170]	77,951	69	-	612,437	Static	General	Wikidata	https://
ConceptNet [171]	28,370,083	50	-	34,074,917	Static	General	Wikipedia,OpenCyc,WordNet	https://
ConceptNet100K [172]	78,527	34	-	100,000	Static	General	ConceptNet	https://
DBpedia50K [173]	49,000	654	-	43,756	Static	General	Wikipedia	https://
FB15K-237 [174]	14,541	237	-	310,116	Static	General	Freebase	https://
Hetionet [175]	47,031	24	-	2,250,197	Static	Biomedical	public datasets	https://
NELL-995 [100]	75,492	200	-	154,213	Static	General	Web	https://
OpenBioLink [176]	180,992	28	-	4,192,002	Static	Biomedical	public datasets	https://
UMLS [177]	135	49	-	6752	Static	Biomedical	human experts	https://
Nation [177]	14	55	-	1,592	Static	Social Sciences	human experts	https://
WN18 [74]	40,943	18	-	151,442	Static	Semantics	WordNet	https://
YAGO3-10 [178]	123,182	37	-	1,079,040	Static	General	YAGO	https://
IMDB-13-3SP [179]	3,244,455	14	30	627,096	Temporal	Movie	IMDB	https://
IMDB-30SP [179]	243,148	14	3	7,923,773	Temporal	Movie	IMDB	https://
YAGO-3SP [179]	27,009	37	3	130,757	Temporal	General	YAGO	https://
DBpedia-3SP [179]	66,967	968	3	201,089	Temporal	General	DBpedia	https://
YAGO11k [134]	10,623	10	189	161,540	Temporal	General	YAGO	https://
Wikidata12k [134]	12,554	24	232	3,419,607	Temporal	General	Wikidata	https://
GDELT-small [180]	500	20	366	3,419,607	Temporal	Social Science	GDELT	https://
ICEWS14 [180]	7,128	230	365	90,730	Temporal	Social Science	ICEWS	https://
ICEWS05-15 [180]	10,488	251	4017	479,329	Temporal	Social Science	ICEWS	https://
ICEWS14 [180]	7,128	230	365	90,730	Temporal	Social Science	ICEWS	https://
FB-IMG [154]	11,757	1231	-	350,293	KG+TXT+IMG	General	Freebase	https://
IMGpedia [181]	14,765,300	44,295,900	-	3,119,207,705	KG+TXT+IMG	General	DBPedia,Wikimedia	https://
MMKG-DB15K [182]	14,777	279	-	99,028	KG+Numeric+IMG	General	Freebase,DBpedia	https://
MMKG-Yago15k [182]	15,283	32	-	122,886	KG+Numeric+IMG	General	Freebase,YAGO	https://
MKG-W [183]	15,000	169	-	42,746	KG+TXT+IMG	General	Wikipedia	https://
MKG-Y [183]	15,000	28	-	26,638	KG+TXT+IMG	General	YAGO	https://
RichPedia [184]	29,985	3	-	119,669,570	KG+IMG	General	WikiPedia	https://
FB15k-237-IMG [161]	14,541	237	-	310,116	KG+IMG	General	Freebase	https://
WN18-IMG [161]	14,541	18	-	151,442	KG+IMG	General	WordNet	https://

Table 1 Dataset statistics of reasoning based on symbolic knowledge bases.

Task definition: Given a knowledge-intensive question, which requires deep understanding and reasoning capability to answer correctly, this type of method is required to leverage the knowledge encoded in parametric knowledge bases to reason for the final answer.

Parametric knowledge bases mainly refer to PLMs. Based on the size of their parameter scales, PLMs can be divided into small language models (SLMs) and large language models (LLMs). SLMs-based reasoning methods [185] need to fine-tune on task-oriented or domain-specific data to enhance their performance. For example, LoP [186] trains RoBERTa [42] on both implicit pre-trained knowledge and explicit free-text statements to symbolic reasoning. The method [67] fine-tunes the GPT-2 [187] to generate full step-by-step solutions to arithmetic reasoning. Though these approaches have shown better performance than traditional rule-based [188, 189], symbolic-based [190, 191], and statistical-based [192, 193] methods, their reasoning ability has been limited by the size of SLMs. Moreover, fine-tuning SLMs requires high-quality training data, which is quite labor-intensive.

Recently, LLMs-based reasoning methods have shown impressive abilities, and prompting is the primary way to interact with LLMs. Compared with SLMs, LLMs possess strong generalization and in-context learning capabilities by providing a few demonstrations (i.e., few-shot learning) or instruction to solve new problems without any demonstrations (i.e., zero-shot learning). Therefore, we mainly investigate reasoning methods based on LLMs.

Many prompting methods have been proposed for reasoning problems [197, 208]. The first attempt is made by [202], which developed a zero-shot prompting method by adding a natural language description of the task in the prompt. Some follow-up methods try to optimize the prompting strategy from the perspective of Chain-of-Thought (CoT) [61], Iterative feedback [62], Problem decomposition [63], Assem-

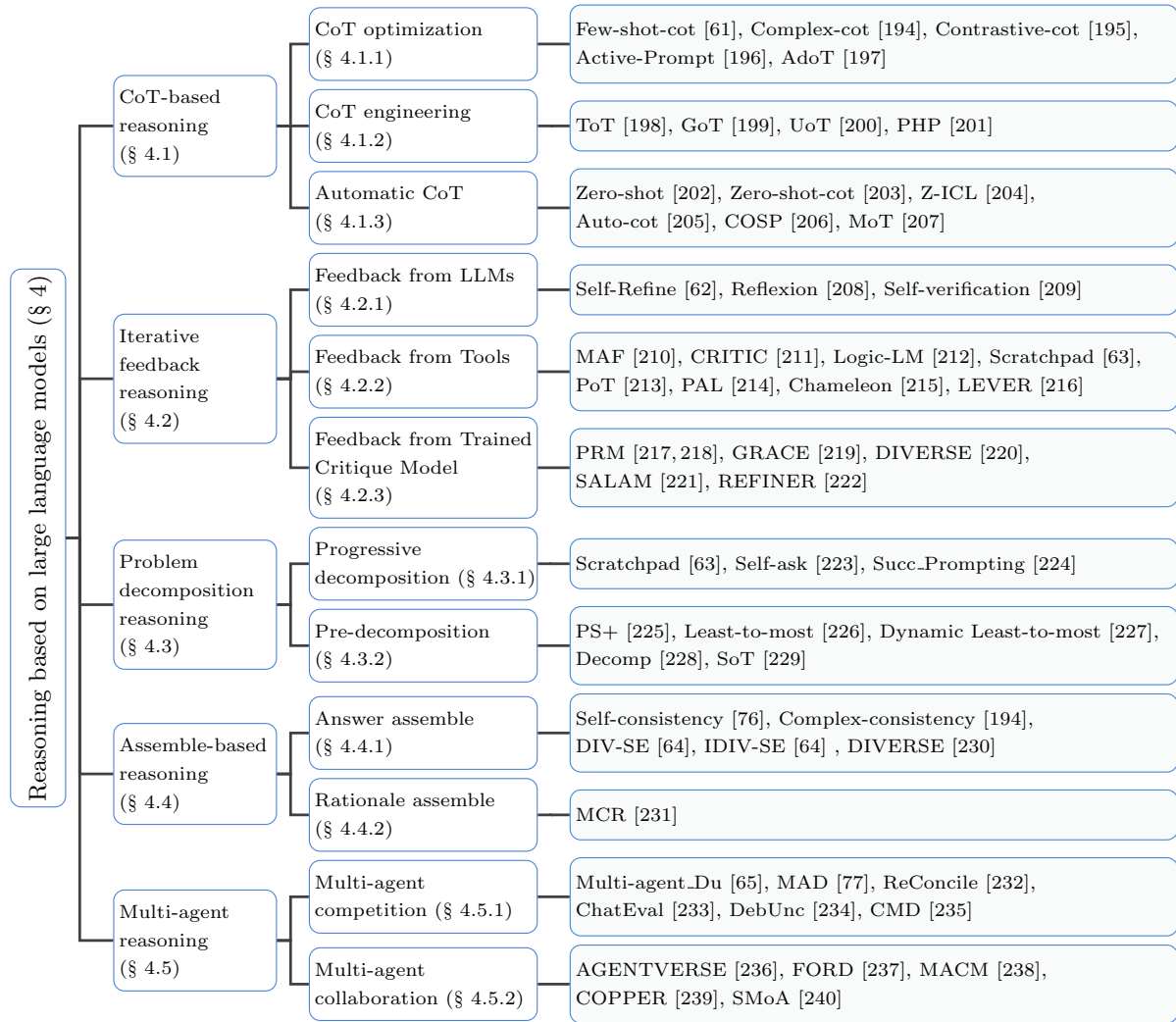


Figure 5 Taxonomy of reasoning methods based on large language models.

ble [76], and multi-agent [77]. These methods are widely applied to mathematical, commonsense, and symbolic reasoning tasks. The overall taxonomy of reasoning methods based on parametric knowledge bases is shown in Figure 5.

4.1 CoT-based reasoning

Recent studies [61, 203, 226] find that generating a series of intermediate reasoning steps (also known as CoT and rationale) significantly improves the ability of LLMs to perform complex reasoning. The intermediate reasoning steps of CoT-series methods contribute to enhancing the logical consistency of the reasoning processes before reaching a conclusion. In this way, CoT-series methods significantly improve the LLMs' ability to perform tasks that require multi-step reasoning and deep understanding. There are three types of CoT-series methods: CoT optimization, CoT engineering, and automatic CoT methods.

4.1.1 CoT optimization

CoT optimization methods construct the question-rationale pairs as demonstrations to guide the LLMs to reason step by step. Few-shot-cot [61] adopts some questions and manually constructs CoT as demonstrations for the first time. Following this line, the CoT optimization methods try to optimize the CoT in demonstrations from different perspectives. For instance, Complex-cot [194] finds that demonstrations with higher reasoning complexity achieve substantially better performance on multi-step reasoning. Hence, it constructs CoT with more reasoning steps in demonstrations. Inspired by how humans can

learn from both positive and negative examples, Contrastive-cot [195] provides both positive and negative demonstrations to enhance the reasoning of LLMs. To determine which questions are the most important and helpful to annotate from a pool of task-specific questions, Active-Prompt [196] proposes an uncertainty-based annotation strategy, which can reduce the model’s uncertainty and help elicit the reasoning ability of LLMs. To solve the mismatch between the question difficulty and the methods’ complexity, AdoT [197] first presents a difficulty measuring approach for questions that computes the syntactic and semantic complexity of their rationales. Then, it proposes a demonstration set construction and a difficulty-adapted retrieval strategy to adaptively construct reasonable demonstrations based on the difficulty of the questions.

4.1.2 *CoT engineering*

CoT engineering methods are devoted to designing more complex CoT reasoning procedures on the top single reasoning process to help LLMs generate more accurate final answers. Inspired by cognitive science, which characterizes problem-solving as a search through a combinatorial problem space, ToT [198] actively maintains a tree of thoughts, where each thought is a coherent language sequence that serves as an intermediate step toward problem-solving. GoT [199] models the information generated by LLMs as an arbitrary graph, where information units are vertices, and edges correspond to dependencies between these vertices. GoT enables combining arbitrary LLM thoughts into synergistic outcomes, distilling the essence of whole networks of thoughts, or enhancing thoughts using feedback loops. The information needed to solve the task is not initially given in many reasoning-related applications. To enhance LLMs in actively seeking information, UoT [200] incentivizes a model to seek information in a way that maximally reduces the amount of information it does not know. To optimize the generated answer progressively, PHP [201] performs automatic multiple interactions between queries and LLMs by using previously generated answers as hints.

4.1.3 *Automatic CoT*

The aforementioned two types of methods achieve excellent performance, but they rely on manually constructed demonstrations, which may generalize poorly between data from different domains. Hence, the automatic CoT methods try to design general instructions to trigger multi-step reasoning or construct pseudo-demonstrations to guide LLMs under a zero-shot setting. For instance, Zero-shot-cot [203] concatenates a simple but effective instruction “Let’s think step by step” after question, which can activate the inherent multi-step reasoning capability of LLMs. To solve the problem that performance drops significantly when no demonstration is available, Z-ICL [204] constructs pseudo-demonstrations from a raw text corpus. It retrieves relevant text from the corpus using the nearest neighbor search and then adjusts the pseudo-demonstrations with physical neighbor and synonym labeling to avoid the copying effect. Auto-cot [205] samples questions with diversity and automatically generates rationales to construct demonstrations. COSP [206] constructs demonstrations from the LLM zero-shot outputs via carefully designed criteria that combine consistency, diversity, and repetition. MoT [207] pre-thinks on the unlabeled dataset and saves the high-confidence thoughts through answer entropy as external memory. During inference, MoT lets the LLM recall relevant memory to help itself reason and answer it.

4.2 Iterative feedback reasoning

Inspired by the human behavior of trial, checking errors, and correcting them during reasoning, some researches focus on utilizing iterative feedback to correct mistakes in the reasoning steps to enhance the reasoning capabilities of LLMs [241]. Using iterative feedback to improve reasoning typically involves three steps: 1) reasoning, 2) critique and feedback, and 3) reasoning refinement. The sources of feedback in these methods mainly include LLMs themselves, various tools (including calculators, search engines, logic tools, and code interpreters), and trained critique models.

4.2.1 *Feedback from LLMs*

Some early studies discovered that LLMs can check errors in their reasoning process, known as the ability of “self-reflection”. Inspired by this, Self-Refine [62] explores how to achieve iterative feedback and refinement based on the LLM itself to improve the quality of reasoning. Reflexion [208] employs an iterative process of “Trajectory→Evaluation→Reflection→Next Trajectory” to iteratively enhance the

reasoning process. Self-verification [209] implements “Forward Reasoning” and “Backward Verification” to validate the reasoning process and selects the highest-quality reasoning results based on evaluation scores. However, it should be noted that recent studies [242,243] have suggested that this “self-reflection” approach may be constrained by the model’s inherent reasoning capabilities, potentially hindering improvements in reasoning quality.

4.2.2 *Feedback from tools*

Research has explored integrating various tools into feedback modules to help correct reasoning errors. Different types of reasoning errors necessitate different tools. For example, calculators are often used to provide precise arithmetic results as feedback for arithmetic tasks [210,211], while search engines are employed to verify factual errors [211,215]. For logical reasoning tasks, Logic-LM [212] suggests using logical tools like First-order Logic Provers to identify logic errors and provide feedback. Moreover, considering that the pre-training corpora of LLMs contain a substantial amount of structured code, some studies suggest transforming reasoning tasks into code form and then using code interpreters to provide feedback [63,213–216]. All these methods focus on transforming the output text of LLMs into formats suitable for inputting into various tools, thus obtaining precise external feedback and enhancing the quality of reasoning.

4.2.3 *Feedback from trained critique model*

Earlier studies [217,218] demonstrated that process-based reward models outperform outcome-based reward models when providing feedback for mathematical reasoning tasks. Building on this, GRACE [219] and DIVERSE [220] propose training critic models capable of selecting optimal intermediate reasoning steps as feedback. Furthermore, SALAM [221] and REFINER [222] explore the idea of training models that generate error analyses in natural language form as feedback, which can be used to refine reasoning steps iteratively. All these approaches involve training task-specific critique models, enabling them to fully leverage the available training data for specific reasoning tasks. As a result, they achieve significantly improved feedback quality and corrective effectiveness compared to relying on general-purpose LLMs for feedback.

4.3 **Problem decomposition reasoning**

When solving complex problems, decomposing a problem into multiple simpler or more detailed sub-problems is an important strategy employed by human. The general process of problem decomposition-related methods is decomposing a complex problem into several simpler sub-problems. These sub-problems are then solved one by one. Finally, the answers to the sub-problems are combined to obtain the answers to the original problem. When the task is complex or the individual reasoning steps are hard to learn, this method often yields superior results. However, when dealing with simple problems, further decomposing them may not be very meaningful and increase time overhead. There are two types of problem decomposition methods: progressive decomposition and pre-decomposition methods.

4.3.1 *Progressive decomposition*

Progressive decomposition methods alternately decompose the questions and reason for the sub-questions step-by-step. For instance, Succ_Prompting [224] iteratively decomposes the complex question into the following simple question to answer and then repeats until the complex question is answered. Scratchpad [63] allows the model to produce an arbitrary sequence of intermediate tokens, which it calls a scratchpad, before producing the final answer. For example, on addition problems, the scratchpad contains the intermediate results from a standard long addition algorithm. Self-ask [223] asks itself follow-up questions before answering the initial question, which will narrow the compositionality gap that models can correctly answer all sub-problems but not generate the overall solution.

4.3.2 *Pre-decomposition*

Unlike the progressive decomposition methods, the pre-decomposition methods decompose the questions before reasoning. PS+ [225] devises a plan to divide the entire task into smaller subtasks and then carries out the subtasks according to the plan. Least-to-most [226] breaks down a complex problem into

a series of simpler subproblems and then solves them in sequence. To promote the practicality of Least-to-most, Dynamic Least-to-most [227] obtain the problem reduction via a multi-step syntactic parse of the input. Furthermore, it dynamically selects exemplars from a fixed pool such that they collectively demonstrate as many parts of the decomposition as possible. Decomp [228] argues that few demonstrations of the complex task aren't sufficient for current models to learn to perform all necessary reasoning steps as tasks become more complicated. Hence, it solves complex tasks by instead decomposing them into simpler sub-tasks and delegating these to sub-task specific LLMs, with both the decomposer and the sub-task LLMs having their own few-shot prompts. SoT [229] guides the LLM to derive a skeleton first by itself. Based on the skeleton, the LLMs then complete each point in parallel.

4.4 Assemble-based reasoning

The core idea of assemble-related methods is that a complex reasoning problem typically admits multiple different ways of thinking, leading to its unique correct answer [76]. As shown in Figure 4 (e), typical answer assemble methods first generate multiple different rationales with answers and then choose the most consistent one as the final answer. Furthermore, a few rationale assemble methods try to leverage the difference between multiple reasoning processes to enhance reasoning performance. The assemble-related methods demonstrate excellent performance. Moreover, they can be easily integrated with other classes of methods, such as CoT-series. However, due to the need to generate multiple reasoning processes, the overhead of this class of methods is relatively high.

4.4.1 Answer assemble

Self-consistency [76] first samples a diverse set of rationales instead of only taking the greedy one and then selects the most consistent answer by marginalizing out the sampled reasoning paths. Instead of voting among all rationales, Complex-consistency [194] votes among top-K complex rationales with more steps. To leverage variations of the input prompt to introduce the diversity needed for assembling, DIV-SE [64] automatically improves prompt diversity by soliciting feedback from the LLM to ideate approaches that are apt for the problem. Then, it assembles the diverse prompts across multiple inference calls. To reduce the inference costs of DIV-SE, IDIV-SE [64] combines all approaches within the same prompt and aggregates all resulting outputs to leverage diversity. Similar to DIV-SE [64], DIVERSE [230] proposes to increase the diversity of rationales by sampling from a single prompt and varying the prompt itself. It first uses a verifier to score the quality of each rationale and guide the voting mechanism. Then, it assigns a fine-grained label to each step of the reasoning path and uses a step-aware verifier to attribute the correctness or wrongness of the final answer to each step.

4.4.2 Rationale assemble

The answer assemble methods generate multiple rationales and aggregate them through a voting mechanism over the final answers, which ignore the information in intermediate steps. Furthermore, although the answer assemble methods perform well, they do not provide a unified explanation for the predicted answer. Hence, MCR [231] focuses on rationale assemble, which leverages the relations between intermediate steps across multiple rationales. MCR mixes information between multiple relations and selects the most relevant facts to generate an explanation and predict the answer. Unlike answer assemble methods, sampled rationales are used not for their predictions but to collect evidence from multiple rationales. MCR concatenates the intermediate steps from each rationale into a unified context, passed to a meta-reasoner model along with the original question. The meta-reasoner model prompts to meta-reason on multiple rationales and produces a final answer with an explanation. In this way, MCR could combine facts from multiple chains to produce the final answer with an explanation of the answer's validity.

4.5 Multi-agent reasoning

Multi-agent reasoning draws inspiration from *society of minds* concepts [244] found in multi-agent systems. In contrast to single-agent methods, such as CoT and ToT, multi-agent reasoning methods emphasize the diversity of ideas and the importance of communication, adversarial interaction, and collaboration among multiple agents. In the reasoning stage, multi-agents express their individual viewpoints and interact in various ways (such as through debate, collaboration, and community communication) to

arrive at a final solution. The divergent thinking of multi-agent reasoning determines that (i) The distorted thinking of one agent can be rectified by other agents, (ii) the supplementation of one agent’s resistance to change by others, and (iii) the reception of external feedback by each agent from others.

A limitation of multi-agent reasoning is that it requires more time cost, as agents often need to participate in multiple rounds of interaction to present and counter arguments. Additionally, current LLM-based agents may struggle to maintain coherence and relevance in long-context scenarios, leading to potential misunderstandings and context loss. Enhancing the long-text modeling capabilities of large language models remains challenging for future research. Multi-agent reasoning methods can be categorized into multi-agent competition and multi-agent collaboration.

4.5.1 *Multi-agent competition*

Multiple LLM-based agents conduct independent thinking and reasoning in multi-agent competition methods. When agents hold differing opinions, they examine each other’s responses and adapt their answers accordingly. Through several rounds of adversarial interaction, the multi-agent system ultimately reaches a conclusion that satisfies the internal logic of each agent while aligning with the feedback from other agents.

Recently, abundant multi-agent competition methods are designed to explore how to unleash the potential of multi-agent systems. For instance, Multi-agent_Du [65] is a role-symmetric multi-agent competition framework where different agents engage in spontaneous discussion. MAD [77] introduced different roles, such as judges and debaters, into the debate process. This represents a role-asymmetric multi-agent debate architecture. RECONCILE [232] facilitates deeper multi-agent discussions by introducing confidence assessments and persuasive explanations in the form of roundtable meetings. ChatEval [233] adopts three distinct communication strategies within its diversified role communication process: one-on-one, simultaneous-talk, and simultaneous-talk-with-summarizer. DebUnc [234] enhances the reliability of multi-agent debates by quantifying and conveying agents’ uncertainties throughout the debate process, reducing the hallucination phenomena often associated with LLMs. Additionally, research evaluating multi-agent competition has shown that, as demonstrated by CMD [235], effective prompt engineering can enable a single agent to achieve performance comparable to multi-agent discussions. Multi-agent discussions, however, have a distinct advantage in contexts lacking examples, and discussions involving multiple LLMs can enhance the performance of weaker LLMs.

4.5.2 *Multi-agent collaboration*

Multi-agent LLM collaboration involves agents working together cooperatively to solve a given problem. AGENTVERSE [236] simulates the problem-solving process of human groups through mechanisms such as expert recruitment, collaborative decision-making, and tool utilization. FORD [237] explores the issue of mutual consistency among multiple LLMs by introducing a formalized debate mechanism, illuminating both the potential and challenges inherent in LLM collaboration. MACM [238] first abstracts the conditions and objectives of a problem, then employs a multi-agent interaction system to iteratively uncover new conditions that facilitate the achievement of the goals, ultimately solving the problem. COPPER [239] enhances the collaborative capabilities of multi-agent systems based on LLMs through a self-reflection mechanism. This framework involves training a shared reflector and utilizes a counterfactual proximal policy optimization (PPO) mechanism to optimize the quality of reflections. SMoA [240] introduces sparse to optimize the fully connected structures commonly found in traditional multi-agent methods, thereby balancing performance and computational cost. Social_Agent [245] explores the collaboration mechanisms among agents and analyzes these mechanisms from a social psychology perspective.

4.6 Datasets

In this section, we have compiled statistics on some commonly-used datasets related to reasoning based on parametric knowledge bases, including (1) Domain: The domain of knowledge corresponding to the datasets; (2) # Ques.: Question number; (3) Q&A source: The main construction methods of questions and answers. They are mainly divided into three categories: “Generate”, “Expert”, and “Crowdsourcing”. “Generate” refers to the design of programs for automated generation, “Expert” refers to direct crawling from professional websites or carefully designed by domain experts, and “Crowdsourcing” refers to completion by crowdsourcing workers with general cultural levels; (4) Rationale: Whether they contain

rationales; (5) Answer type: The form of answers; (6) Links: The storage address of the datasets. The statistical results are shown in Table 2.

Datasets	Domain	# Ques.	Q&A source	Rationale	Answer type	Links
AQUA [246]	Arithmetic	254	Generate	✓	Option	https://
GSM8K [247]	Arithmetic	1319	Crowdsourcing	×	Number	https://
SVAMP [248]	Arithmetic	1000	Generate	✓	Number	https://
AddSub [249]	Arithmetic	395	Expert	✓	Number	https://
MultiArith [250]	Arithmetic	600	Expert	✓	Number	https://
SingleEq [251]	Arithmetic	508	Expert	✓	Number	https://
Last Letters [61]	Symbolic	500	Generate	×	String	https://
Coin Flip [61]	Symbolic	500	Generate	×	Yes/No	https://
CommonsenseQA [252]	Commonsense	1221	Crowdsourcing	×	Option	https://
StrategyQA [253]	Commonsense	2290	Crowdsourcing	✓	Yes/No	https://

Table 2 Dataset statistics of reasoning based on parametric knowledge bases.

5 Collaborative reasoning based on symbolic and parametric knowledge bases

In this section, we are devoted to investigating collaborative reasoning methods based on symbolic and parametric knowledge bases. These methods mainly perform reasoning in the form of question answering. Generally, given a knowledge-intensive question, this type of method is required to leverage the knowledge in symbolic and parametric knowledge bases collaboratively to reason for the correct answer.

Based on the structure of the symbolic knowledge bases, these question answering tasks can be categorized into graph-based reasoning, table-based reasoning, text-based reasoning, and heterogeneous reasoning. The symbolic knowledge bases store knowledge in the form of structured graphs, structured tables, and unstructured text in graph-based, table-based, and text-based reasoning tasks, respectively. In particular, heterogeneous reasoning investigates how to leverage symbolic knowledge from multiple heterogeneous symbolic knowledge bases, such as KGs, tables, and text. The overall taxonomy of collaborative reasoning methods based on symbolic and parametric knowledge bases is shown in Figure 6.

5.1 Graph-based reasoning

The tasks of graph-based reasoning include knowledge graph question answering (KGQA) and temporal knowledge graph question answering (TKGQA), where symbolic knowledge is stored in SKGs and TKGs, respectively.

5.1.1 Knowledge graph question answering

Task definition: Given a question and an SKG, KGQA methods are required to understand the intent of the question via the parametric knowledge bases and retrieve the entity nodes from the SKG as answers.

KGQA methods can be categorized into two classes: semantic parsing-based (SP-based) methods and information retrieval-based (IR-based) methods [36]. SP-based methods aim to parse the questions into the logical forms (such as SPARQL, S-expression and query graph) to yield the correct answer. IR-based methods construct a question-specific subgraph of the SKG and retrieve the most matching answers. In recent years, most methods utilize PLMs to integrate a substantial amount of external knowledge. Due to the introduction of parametric knowledge, the level of intelligence has been significantly enhanced, leading to considerable improvements in accuracy and task versatility.

5.1.1.1 Semantic parsing-based methods

SP-based methods aim to learn the semantic matching between natural language questions and logical forms, which mainly involves the following steps: the method understands the natural language question, converts it into a logical form, aligns it with the existing knowledge in given SKG, and finally executes it to derive the correct answer. Early SP-based methods only apply to independent and identically distributed scenarios and perform poorly in solving problems that require commonsense knowledge and

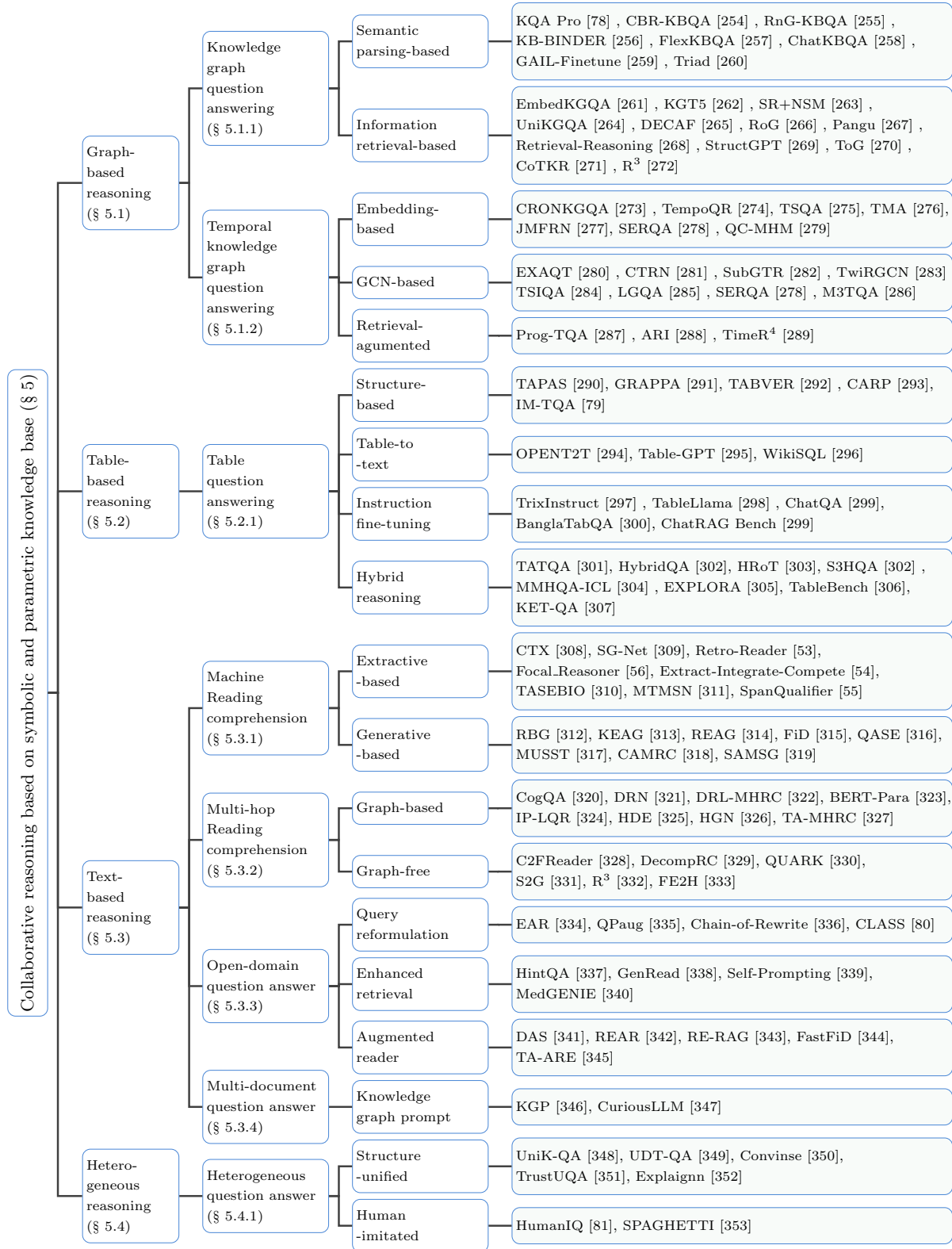


Figure 6 Taxonomy of collaborative reasoning methods based on symbolic and parametric knowledge bases.

exhibit relatively weak intelligence. So most SP-based methods fine-tune PLMs on specific data to convert natural language questions to logical form.

For example, KQA Pro [78] introduces a compositional and interpretable programming language KoPL to represent the reasoning process of complex questions fine-tunes BART [69] to achieve compositional reasoning. CBR-KBQA [254] uses ROBERTA-base [42] to encode each question independently and generate a logical form for a new question by retrieving cases that are relevant to it with the pre-trained ROBERTA-base weights. RnG-KBQA [255] introduces T5 [68] to construct the final logical form based on the questions and the high-ranked candidate logical forms, demonstrating excellent performance even when dealing with questions involving unseen schema items.

Several SP-based methods utilize LLMs to parse natural language questions into logical forms in a few-shot in-context learning setting. For example, KB-BINDER [256] generates drafts with LLM as preliminary logical forms and then binds the entities, relations, and schema items of the drafts to SKG iteratively. FlexKBQA [257] introduces a self-training approach with execution guidance, using the LLM to convert logical forms into natural language questions and utilizing unlabeled user questions iteratively.

In addition, fine-tuning LLM also makes sense. For example, ChatKBQA [258] proposes generating the logical form with fine-tuned LLMs first, then retrieving and replacing entities and relations through an unsupervised retrieval method. GAIL-Finetune [259] fine-tunes Llama-2-7B [354] to produce expert-level sample and evaluate the authenticity and relevance of the sequences to tackle the challenges in low-resource KGQA scenario. Triad [260] utilizes an LLM-based agent with three different roles for KBQA tasks.

In general, knowledge for KGQA tasks not only includes professional knowledge that can be modeled through SKGs and limited training samples but also general knowledge. Therefore, fine-tuning PLM, or directly utilizing the powerful LLM, can overcome difficulties in collecting training samples and understanding diversified natural language questions in the real world. However, there are also deviations between parametric general knowledge and symbolic professional knowledge, and exploration of this issue is still in the early stages.

5.1.1.2 *Information retrieval-based methods*

IR-based methods treat the KGQA task as a binary classification of nodes in the SKG and demonstrate superior performance compared to SP-based methods, although they sacrifice interpretability. Early IR-based methods perform poorly because they lack sufficient understanding of user questions and the guidance of common sense knowledge during the reasoning process. Hence, recent advanced IR-based methods generally apply PLMs to model questions and subgraphs of the SKG.

For example, EmbedKGQA [261] has been proposed to handle SKG sparsity, where ComplEx [91] embeddings are trained to represent SKG elements and RoBERTa [42] embedding is used to represent the question. KGT5 [262] considers both SKG reasoning and KGQA as sequence-to-sequence tasks, where a simple Transformer that has the same architecture as T5-small [68] has been trained to achieve excellent performance. SR+NSM [263] utilizes RoBERTa [42] to encode the question and relations in SKG iteratively to expand paths. Then it could construct subgraph with low size but high answer coverage to find answers. UniKGQA [264] combines a PLM with an ultra-simple GNN to transfer the retrieved knowledge to the reasoning phase. DECAF [265] constructs retriever and reader based on FiD-large retriever [355] to generate both logical forms and direct answers jointly.

Recently, several methods have adopted the retrieval-augmented generation (RAG) paradigm because LLMs can simultaneously model both user questions and SKG elements to perform simple deductive reasoning. For example, RoG [266] proposes a planning-retrieval-reasoning framework, which fine-tunes LLaMA2-Chat-7B [354] with relation paths and valid reasoning paths in SKGs. Then it can generate reasoning paths for faithful reasoning. Similarly, Retrieval-Reasoning [268] decomposes the problem into retrieval and reasoning modules and then fine-tunes LLM at three levels: entity, relation, and graph. Pangu [267] leverages the discriminative capabilities of the LLM for context-based language understanding. The symbolic agent explores SKG to construct effective plans incrementally, and the LLM agent evaluates the reasonableness of candidate plans to guide the search process. StructGPT [269] solves KGQA based on structured data, where the facts in SKG could be linearized into LLM to reason naturally. ToG [270] treats LLM as an agent capable of exploring SKG and performing reasoning with retrieved knowledge. The agent iteratively executes beam search on the KG, discovers the most promising

reasoning paths, and returns the most likely reasoning results. CoTKR [271] rewrites retrieved subgraphs into natural language formats comprehensible to LLMs. R³ [272] surfaces the commonsense knowledge relevant to the question from LLMs and uses it to guide the SKG pruning to find answers.

Overall, IR-based methods require identifying supportive evidence to answer the question directly. In contrast, the parametric knowledge within PLM is used to learn the user’s intent and retrieve the relevant subgraph, reducing semantic noise while maintaining a high recall rate of supportive evidence. More importantly, introducing PLM can make the retrieval and denoising processes more transparent and maintain certain performance in few-shot or even zero-shot scenarios.

5.1.2 Temporal knowledge graph question answering

Task definition: Given a temporal question and a TKG, TKGQA methods are required to understand the intent and temporal constraints of the question via the parametric knowledge bases and retrieve the entity nodes or timestamps from the TKG as answers.

The primary solution for TKGQA involves integrating the knowledge of TKGs and PLMs. Some methods use PLM and pre-trained TKG embeddings to match questions and entities and timestamps, known as embedding-based methods. Some other methods utilize PLM and GCN [167] to learn the features of TKG elements, known as GCN-based methods. In recent years, there have also been methods utilize LLMs by converting the facts from TKGs into text, which are named retrieval-augmented methods because they primarily enhance the reasoning capabilities of LLMs by accurately retrieving useful knowledge from TKGs.

5.1.2.1 Embedding-based methods

Embedding-based methods utilize existing TKG model to derive embeddings of TKG elements and utilize PLM to derive embeddings of questions. By modeling the matching between them, the answers could be inferred based on the similarity between embeddings of given question and all TKG elements.

A classic embedding-based method is CRONKGQA [273], which builds on EmbedKGQA [261] by employing an advanced temporal knowledge graph reasoning model to derive embeddings of entities and timestamps. It uses PLM to obtain embeddings of entity/time mentioned in the question, models the question as a “virtual relation”, and predicts missing entities and timestamps in the TKG. Similarly, QC-MHM [279] is more refined in handling questions and TKGs. It first injects temporal order information into timestamp embeddings, modifying TComplex [131] to obtain the embeddings of the entity, relation, and timestamp. It then inputs the sentence and SPO (subject, predicate or relation, and object) into Sentence-BERT [356] to obtain the embedding vectors and model the matching between questions and SPO. Other similar methods include TempoQR [274], TSQA [275], TMA [276], JMFRN [277], and SERQA [278].

The advantage of embedding-based methods lies in their more natural thought process, with a lighter model that can better model the question and match it with TKG. However, the disadvantage is that they struggle with complex temporal constraints and temporal relational terms. Complex temporal constraints require the utilization of multiple factual quadruples, and temporal relational terms are highly sensitive to model.

5.1.2.2 GCN-based methods

GCN-based methods do not leverage existing TKG reasoning models to incorporate knowledge from TKGs, and they learn the matching between questions and TKG subgraphs. They utilize PLM and GCN [167] to learn embeddings of questions and features of a subgraph of TKG, respectively. Then, they cast answer prediction into a node classification task. Compared to embedding-based methods, GCN-based methods can model various constraints within complex problems and retain a more complete subgraph of the TKG.

For example, EXAQT [280] utilizes fine-tuned BERT models and GCN [167] to identify relevant facts. It specifically employs Group Steiner Trees to compute question-relevant compact subgraphs within the KG. Additionally, relational graph convolutional network has been constructed to predict answers. Similarly, GenTKGQA [357] first leverages a LLM and a pre-trained temporal graph neural network to model question and extract information from the subgraph, respectively. Then, it performs instruction tuning to enable complex temporal reasoning. Other similar methods include SubGTR [282], TSIQA [284], CTRN [281], TwiRGCN [283], M3TQA [286] and LGQA [285].

5.1.2.3 Retrieval-augmented methods

Retrieval-augmented methods primarily rely on the in-context learning ability and powerful generation capability of LLM, where there is no need for instruction tuning. At the same time, the most relevant knowledge is retrieved from the TKG and combined with the LLM’s knowledge to complete complex reasoning tasks collaboratively. For example, Prog-TQA [287] uses a LLM to understand questions and generate corresponding program drafts with symbolic operators as logical forms, given a few examples. With a self-improvement strategy, the quality of these logical forms is enhanced to yield the final answers. Similarly, ARI [288] improves the LLM’s capacity to integrate abstract methodologies derived from historical experience. Timer⁴ [289] differentiates the knowledge of TKG into temporal knowledge and factual knowledge and then improves retrieval accuracy through modules such as retrieval, rewriting, and ranking.

The advantage of retrieval-augmented methods lies in their emphasis on retrieving TKG elements. By textualizing the knowledge within the TKG to reduce hallucinations in LLMs and enhance their reasoning capabilities, they achieve breakthroughs in performance than embedding-based methods and GCN-based methods. However, due to LLM’s insufficient sensitivity to temporal knowledge, they still face challenges regarding knowledge integration and complex temporal reasoning.

5.1.3 Datasets

In this section, we have compiled statistics on some commonly-used datasets related to graph-based reasoning, including (1) Domain: The domain of knowledge corresponding to the datasets; (2) Question type: The type of questions in the datasets; (3) # Ques.: Question number; (4) Q&A source: The main construction methods of questions and answers. They are mainly divided into three categories: “Generate”, “Expert”, and “Crowdsourcing”. “Generate” refers to the design of programs for automated generation, “Expert” refers to direct crawling from professional websites or carefully designed by domain experts, and “Crowdsourcing” refers to completion by crowdsourcing workers with general cultural levels; (5) KG: Specific KG (SKG or TKG) they use; (6) Links: The storage address of the datasets. The statistical results are shown in Table 3.

Datasets	Domain	Question type	# Ques.	Q&A source	KG	Links
Free917 [358]	General	Static	917	Expert	Freebase	https://
WebQuestions [359]	General	Static	5,810	Crowdsourcing	Freebase	https://
WebQuestionsSP [360]	General	Static	4,737	Crowdsourcing	Freebase	https://
ComplexQuestions [361]	General	Static	2,100	Expert	Freebase	https://
MetaQA/1-hop [362]	Movie	Static	116,045	Generate	Wikipedia	https://
MetaQA/2-hop [362]	Movie	Static	148,724	Generate	Wikipedia	https://
MetaQA/3-hop [362]	Movie	Static	142,744	Generate	Wikipedia	https://
QALD [363]	General	Static	806	Expert	DBpedia	https://
LC-QuAD [364]	General	Static	5,000	Generate	DBpedia	https://
LC-QuAD2.0 [365]	General	Static	5,000	Generate	DBpedia	https://
TempQuestions [366]	General	Temporal	1,271	Expert	Freebase	https://
TimeQuestions [280]	General	Temporal	16,181	Expert	Wikidata	https://
CRONQUESTIONS [273]	General	Temporal	410,000	Expert	Wikidata	https://
Complex-CRONQUESTIONS [282]	General	Temporal	45,821	Expert	Wikidata	https://
MultiTQ [367]	Social Science	Temporal	500,000	Expert	ICEWS	https://

Table 3 Dataset statistics of graph-based reasoning.

5.2 Table-based reasoning

The task of table-based reasoning mainly refers to table question answering (Table QA), where symbolic knowledge is stored in the structured table.

5.2.1 Table question answering

Task definition: Given a question and a table, Table QA methods are required to understand the intent of the question via the parametric knowledge bases and find the correct answer from the table. The table QA methods can be divided into: structure-based, table-to-text, instruction fine-tuning, and hybrid reasoning methods.

5.2.1.1 *Structure-based methods*

Structure-based methods parse the table structure and utilize the relationships between rows, columns, and cells to find the final answers. This approach has proven effective in extracting structured insights by focusing on the table’s inherent organization.

Many methods primarily focus on understanding the semantic association between table headers and data areas. For example, Müller [368] proposes encoding a table into a graph and employing GNNs and pointer networks to select answers and address sequential questions directly from the table. Similarly, TAPAS [290] flattens tables, encodes structure through multiple positional embeddings, and uses pre-training on text-table pairs to predict cell selections and aggregation operators for Table QA. GRAPPA [291] takes a different approach by inducing synchronous context-free grammar to generate synthetic data, combining pre-training on masked language modeling and SQL semantic prediction. This enables GRAPPA to utilize structural knowledge effectively for semantic parsing, converting user inputs into executable programs, and enhancing performance across various datasets.

However, while these methods excel in querying structured data accurately, they face significant limitations when dealing with complex tables, such as nested structures, multi-table setups, or untitled data. They also struggle with flexible questions or generative tasks, highlighting a gap in adaptability. To address these challenges, a few methods have been introduced. For example, TABVER [292] integrates arithmetic reasoning with natural logic reasoning systems, enabling tabular fact-checking tasks to overcome the limitations of symbolic reasoning models and natural logic systems in handling arithmetic operations. Similarly, CARP [293] employs mixed-modal reasoning chains to explicitly model intermediate reasoning steps, improving the interpretability of the model’s reasoning process. Despite these advancements, studies reveal that when rows and columns are rearranged to create new examples, the performance of large models declines significantly. This suggests that current approaches to table structure understanding lack robustness and fail to adapt effectively to structural changes [369].

5.2.1.2 *Table-to-text methods*

Table-to-text methods first convert tabular data into natural language descriptions. Then, it leverage the natural language descriptions to support the reasoning process. Table-to-text methods is flexible and capable of handling complex queries. However, its performance heavily depends on training quality and demands significant computational resources while offering limited explainability.

To address these challenges, several models and frameworks have been proposed. For instance, OPENT2T [294], an open-source toolkit, facilitates the replication and comparison of existing systems, driving innovation in developing new models. Another example is Table-GPT [295], which proposes a “table-tuning” paradigm to improve language models’ performance on table-related tasks. Additionally, methods for summarizing table contents enable users to complete question-answering tasks without needing to browse individual entries in the table [370–372]. Despite these advancements, current research primarily focuses on surface-level achievements, with limited attention to logical reasoning. While existing methods address issues of surface authenticity, they often restate data facts without demonstrating robust reasoning or generalization capabilities. Logical natural language generation [373] aims to bridge this gap by enabling models to generate logically inferred natural language statements from facts in open-domain, structured tables.

5.2.1.3 *Instruction fine-tuning methods*

Recently, efforts have been made to enhance LLMs’ ability to process table data by developing specialized instruction fine-tuning datasets. Notable among these are TrixInstruct [297] and TableLlama [298], which utilize datasets that cover diverse, realistic tables and related tasks. After fine-tuning on these datasets, LLMs show significant improvements in handling table-based questions. Additionally, ChatQA [299] employs a two-stage instruction fine-tuning strategy, yielding substantial gains in table-related tasks. The first stage involves supervised fine-tuning on diverse instruction datasets, while the second stage, context-enhanced instruction fine-tuning, incorporates table QA and other high-quality QA datasets to further refine the model’s conversational QA capabilities in context-specific scenarios.

5.2.1.4 Hybrid reasoning methods

Hybrid reasoning methods combine the above methods with more advanced deep learning technology to process table QA, making it well-suited for multi-step reasoning and complex tasks [374, 375].

Building on this foundation, HRoT [303] introduces retrieval thinking and LLM as a retrieval module, which reconstructs the tables and constructs prompts. In the reasoning stage, it guides the model to retrieve evidence in texts and tables gradually, avoiding the use of irrelevant information and significantly improving the effectiveness of the retrieval process. S3HQA [302] proposes a three-stage framework—comprising retrieval, selector, and reasoner—where LLMs are employed as generative components in the reasoning stage. Expanding on this concept, MMHQA-ICL [304] converts tables and images into text and sends it to the Classifier and Retriever Module for the question type and retrieved documents. Then the Prompt Generator Module builds an LLMs input using the questions and retrieved data. Finally, LLMs give the exact answer. This end-to-end LLM prompt method performs better than the baseline method on the Multimodal QA dataset. While the hybrid reasoning method is effective in solving complex table question answering tasks, it comes with drawbacks such as being complex to implement, requiring high computing resources, and in certain situations, needing manual rule design or model fine-tuning, EXPLORA [305] offers a distinct perspective which can save resources by reducing the number of LLM calls. As a static subset selection method, it uses a scoring function to select examples, bypassing reliance on LLM parameters or outputs.

5.2.2 Datasets

In this section, we have compiled statistics on some commonly-used datasets related to table-based reasoning, including (1) Domain: The domain of knowledge corresponding to the datasets; (2) # Ques.: Question number; (3) Q&A source: The main construction methods of questions and answers. They are mainly divided into three categories: “Generate”, “Expert”, and “Crowdsourcing”. “Generate” refers to the design of programs for automated generation, “Expert” refers to direct crawling from professional websites or carefully designed by domain experts, and “Crowdsourcing” refers to completion by crowdsourcing workers with general cultural levels; (4) Table source: The source of tables; (5) Links: The storage address of the datasets. The statistical results are shown in Table 4.

Datasets	Domain	# Ques.	Q&A source	Table source	Links
HiTab [376]	General	14 K	statistical reports, Wikipedia	Crowdsourcing	https://
FeTaQA [377]	General	10 K	Crowdsourcing	Wikipedia	https://
WikiSQL [296]	General	24 K	Expert	Wikipedia	https://
TableInstruct [298]	General	1.24 M	Expert	Wikipedia, statistical scientific reports	https://
FEVEROUS [378]	General	87 K	Expert	Wikipedia	https://
TableBench [306]	General	0.8 K	Crowdsourcing	existing datasets	https://
BanglaTabQA [300]	General	19 K	Generate	Wikipedia	https://
ChatRAG Bench [299]	General	29.2 K	Expert	Internet	https://
KET-QA [307]	General	13 K	Crowdsourcing	Wikidata	https://
IM-TQA [79]	General	1.2 K	Crowdsourcing	published studies	https://

Table 4 Dataset statistics of table-based reasoning.

5.3 Text-based reasoning

The tasks of text-based reasoning include Machine reading comprehension (MRC), Multi-hop reading comprehension (MHRC), also known as multi-hop question answering), Open-domain question answering (ODQA), and Multi-document question answering (MDQA), where symbolic knowledge is stored in unstructured text.

5.3.1 Machine reading comprehension

Task definition: Machine reading comprehension aims to evaluate a machine’s ability to understand language effectively. MRC methods are presented with one or more text passages and are then required to answer questions based on the provided passages [53, 309, 379, 380]. As highlighted by [56, 381], improving MRC performance involves developing several skills, including numerical reasoning, commonsense reasoning, and logical reasoning.

The research on machine reading comprehension has garnered significant interest over the past decade. The MRC tasks have progressed from the initial cloze-style tests [382,383] to span-based answer extraction from passages [384,385], as well as to multiple-choice [386] and free answering formats [387]. In the early years, rule-based methods [388,389] focus on designing heuristic algorithms. These algorithms are specifically tailored to the grammar of a language and are used to assist in learning, discovery, or problem-solving. By employing trial-and-error techniques, they aim to find evidence within a given sentence to answer a question. Statistic-based methods [390] involves quantifying occurrences of words and utilizing these numerical representations to infer potential answers. Recently, numerous studies [391,392] focus on leveraging machine learning methods to extract the features in the questions and passages, which enhance the machine’s ability to understand and process text automatically.

With the development of deep learning, various attention-based methods [393–396] are proposed to facilitate interactions between passages and questions. Recently, PLMs have achieved significant success in MRC tasks. The PLMs-based models demonstrate a strong ability to capture contextual and sentence-level language representations, which notably improve the benchmark performance of current MRC systems [53,56,309]. In line with this trend, our focus is primarily on PLMs-based MRC methods, which combine symbolic knowledge in text passages and parameter knowledge in PLMs. The PLMs-based MRC methods can be further divided into Extractive MRC and Generative MRC [316]. The Extractive MRC methods try to predict the start and end positions of answers directly from the context, while the Generative MRC methods are devoted to generating answers by reformulating information across the context.

5.3.1.1 *Extractive MRC methods*

Recent MRC research predominantly focuses on extractive question answering using encoder-only PLM, which predicts the start and end positions of answers directly from the context. For instance, CTX [308] adopt PLMs to independently obtain paragraph representations conditioned with the current question, previous questions, and previous answers. To extract the answer span, the start and end positions of the current answer are predicted based on the concatenation of the paragraph representations encoded in the previous step.

In this category of methods, PLMs are mainly used as encoders to extract general features from text paragraphs and questions. The focus is on designing downstream models to extract task-oriented features. For instance, SG-Net [309] leverages syntactic guidance in text modeling, achieving substantial performance gains in MRC by introducing explicit syntactic constraints in the attention mechanism. Inspired by how humans solve reading comprehension questions, Retro-Reader [53] integrates two reading and verification strategies stages. First, it uses sketchy reading to quickly grasp the relationship between the passage and the question, forming an initial judgment. Then, it conducts intensive reading to verify the answer and provide the final prediction. Focal_Reasoner [56] extracts fact units from raw texts via syntactic processing and constructs a supergraph. Then, it performs reasoning over the supergraph and a logical fact regularization and aggregates the learned representation to decode the correct answer. Extract-Integrate-Compete [54] iteratively selects complementary evidence with a novel query updating mechanism and adaptively distills supportive evidence, followed by a pairwise competition to push models to learn the subtle difference among similar text pieces.

In the methods mentioned above, the system is mainly expected to extract a single answer from the passage for a given question. However, in many scenarios, questions may have multiple answers scattered in the passages, and all the answers should be found to answer the questions completely. For extracting answers with multi-span, TASEBIO [310] transfers MRC to a sequence tagging task, predicting whether each token is part of the answer. MTMSN [311] combines a multi-type answer predictor designed to support various answer types (e.g., span, count, negation, and arithmetic expression) with a multi-span extraction method for dynamically producing one or multiple text spans. SpanQualifier [55] presents a novel span-centric scheme to generate representations for all spans in the context and predicts a qualification threshold. Furthermore, it designs a global loss function to jointly optimize overall spans instead of independently optimizing loss on each individual span, which avoids the influence of label imbalance on training the proposed span-centric scheme.

5.3.1.2 *Generative MRC methods*

Recently, significant progress has been made in controllable text generation. Beyond extractive meth-

ods, there is also growing interest in applying generative language models for extractive MRC, which generate answers by reformulating information across the context. For instance, RBG [312] combines a Seq2Seq language model-based generator with a machine reading comprehension module. The reader produces an evidence probability score for each sentence, which will be integrated with the generator for final distribution prediction. KEAG [313] composes a natural answer by exploiting and aggregating evidence from all four information sources available: question, passage, vocabulary, and knowledge. During the process of answer generation, KEAG adaptively determines when to utilize symbolic knowledge and which fact from the knowledge is useful. REAG [314] incorporates an extractive mechanism into a generative model to leverage relevant information to a given question in the contextual passage. Specifically, REAG adds an extraction task on the encoder to obtain the rationale for an answer, which is the most relevant piece of text in an input document to the given question. FiD [315] adopts a simple attention-based inference strategy to extract answer spans from a seq2seq Transformer model without introducing any additional parameters. It then proposes a joint training strategy by combining the normal generative loss and a span extractive loss by enforcing cross-attention to align with answer span positions within the context passages. QASE [316] proposes a novel adaptation of controlled text generation tailored to the specific challenges of MRC, focusing on the precision and relevance of generated answers. Unlike methods that modify the overall generative process through complex architectural alterations or additional learning mechanisms, QASE directly utilizes the question and context to guide inferences.

For multi-span answers, MUSST [317] combines the benefits of span extraction and the simplicity of a multi-span approach to generate free-form answers. It also provides a comprehensive framework for multi-passage generative MRC, which consists of a passage ranker, a multi-span answer annotator, and a question-answering module. CAMRC [318] proposes an answer making-up method from extracted multi-spans learned as highly confident n-gram candidates in the given passage. Unlike the studies that mainly focus on introducing generative mechanisms, SAMSG [319] focuses on handling the writing form of the answer and proposes a novel non-generative decoder to exploit the results from the extractive decoder fully. It learns to score every word in the given passage for how likely they are in the expected answer, then calculates the score of a candidate span from the words' scores.

5.3.2 Multi-hop reading comprehension

Task definition: Multi-hop reading comprehension methods focus on integrating and reasoning over multiple pieces of evidence to answer complex questions. Unlike single-hop MRC, where questions are typically straightforward, and answers can be derived from one or a few nearby sentences, MHRC involves reasoning chains that traverse multiple sentences or even passages. This requires a deep text understanding and reasoning capability, making it more akin to real-world scenarios.

The key challenge of MHRC lies in its demand for multi-step reasoning, where a model must identify and connect intermediate information to form a coherent reasoning path. This reasoning chain culminates in the extraction of the correct answer, supported by a series of evidence sentences. Therefore, MHRC not only tests a model's ability to find the answer but also its capacity to justify the reasoning process with a clear rationale, presenting evidence as proof of the multi-hop reasoning process. Datasets like HotpotQA [397] have been specifically designed to evaluate such multi-hop reasoning capabilities. They include tasks that expect the model to extract and present evidence sentences, thereby showing a clear reasoning trail. Consequently, MHRC better aligns with real-world scenarios where information is dispersed across long passages or multiple documents, necessitating more comprehensive models than those used for single-hop MRC tasks. In MHRC, the given passages are considered as symbolic knowledge bases, and the parametric knowledge bases (PLMs) are adopted to encode the questions and the passages. In MHRC, methods are primarily divided into two categories: graph-based and graph-free methods [327, 331, 332]. Graph-based methods must construct and reason over a graph structure created from the input data. However, graph-free methods do not rely on such structures and often utilize techniques like self-attention for reasoning.

5.3.2.1 Graph-based MHRC methods

The main idea behind graph-based approaches is to represent the input data, which includes the context and questions, as a graph structure. This involves creating nodes representing entities or pieces of text and edges to denote relationships or co-occurrences between these nodes. The reasoning process in MHRC is then performed through message passing over this graph, which allows the model to simulate

the multi-hop reasoning process as it navigates through various interconnected nodes.

Most current research relies on entity graphs where nodes are formed from the context and question entities. For instance, Entity-GCN [398] compiles scattered information from one or more documents by building an entity graph. In this structure, nodes represent entity mentions, while edges illustrate relationships between these mentions within and across multiple documents. BAG [399] transforms documents into a graph in which nodes are entities and edges are relationships between them. The graph is then imported into graph convolutional networks to learn relation-aware representations of nodes. Furthermore, BAG introduces bi-directional attention between the graph and a query with multi-level features to derive the mutual information for the final prediction. Many studies assume that all contexts are pertinent and ignore the negative impact of irrelevant contexts. To filter out unrelated context, DFGN [400] designs a paragraph-selection module to eliminate unrelated paragraphs. It dynamically builds an entity graph from the question entities to locate relevant supporting entities and text spans. Based on DFGN, DFGN_Dual [401] introduces dual reasoning channels to predict the final answer and supporting facts, respectively, which gain better step-by-step reasoning compared to a single-channel approach. Similarly, SAE [402] incorporates a paragraph-selection step to filter out irrelevant context segments, thereby shrinking the problem space while utilizing sentences as graph nodes.

To further improve the performance, some methods try to emulate the human brain's cognitive processes for multi-hop MRC. For example, inspired by the dual process theory of human [403–406], CogQA [320] builds a cognitive graph in an iterative process by coordinating an implicit extraction module and an explicit reasoning module. The extraction module extracts question-relevant entities to construct the cognitive graph. Then, the reasoning module conducts the reasoning procedure over the graph and collects clues to guide the extraction module in extracting next-hop entities better. DRN [321] designs a query reshaping mechanism that visits a query repeatedly to mimic people's reading habits. It dynamically reasons over an entity graph with graph attention and the query reshaping mechanism to promote its comprehension and reasoning ability.

In addition to the work mentioned above, other approaches design models from different perspectives. For instance, DRL-MHRC [322] proposes an RL-based method capable of learning sequential reasoning across extensive collections of documents to pass a query-aware, fixed-size context subset to existing models for answer extraction. BERT-Para [323] first extracts a discrete reasoning chain over the text, which consists of a series of sentences leading to the answer. It then feeds the extracted chains to a BERT-based QA model to predict the final answer. IP-LQR [324] incorporates phrases in the latent query reformulation to improve the cognitive ability of the proposed method for MHRC.

Some studies try to construct more intricate graphs using multiple node types to encompass the available contextual information in the graph constructions fully. For instance, HDE [325] proposes a heterogeneous document-entity graph, which contains different granularity levels of information, including candidates, documents, and entities in specific document contexts. To aggregate clues from scattered texts across multiple paragraphs, HGN [326] creates a hierarchical graph by constructing nodes on different granularity levels, including questions, paragraphs, sentences, and entities. Furthermore, TA-MHRC [327] uses more helpful information about the context, such as the topic of sentences, the topic of relationships, and the importance and strength of relationships, when filtering paragraphs and constructing the graph. Thus, the proposed graph is a weighted graph with four types of nodes and six types of edges to cover the complete information of the context.

5.3.2.2 Graph-free MHRC methods

Compared with the graph-based methods, graph-free methods avoid the explicit construction of graph structures. Instead, they typically rely on more straightforward architectures, possibly leveraging PLMs and self-attention mechanisms to process input data. It is worth noting that C2FReader [328] finds that graph structure can play an important role only when the PLMs are used in a feature-based manner. While the PLMs are used in the fine-tuning approach, the graph structure may not be helpful.

Although graph-free approaches suffer from a performance gap compared to the best graph-based models, numerous methods try to merge the gap by designing powerful mechanisms. Decomprc [329] first decomposes the multi-hop question into several single-hop sub-questions according to a few reasoning types in parallel. Then, Decomprc leverages a single-hop reading comprehension model for every reasoning type to answer each sub-question and combines the answers according to the reasoning type. Finally, Decomprc leverages a decomposition scorer to judge which decomposition is the most suitable

and outputs the answer from that decomposition as the final answer. QUARK [330] scores individual sentences from an input set of paragraphs based on their relevance to the question. Then, it feeds the highest-scoring sentences to a span prediction model to produce an answer to the question. Finally, it scores sentences from the input set of paragraphs again to identify the supporting sentences using the answer. S2G [331] retrieves evidence paragraphs in a coarse-to-fine manner, incorporated with two novel attention mechanisms to restrict the receptive fields of each token according to the nature of each specific task. Inspired by the F1 score, R³ [332] develops an F1 Smoothing mechanism to calculate the significance of each token within the smooth distribution. Furthermore, it incorporates curriculum learning [407] and devises the linear decay label smoothing algorithm, gradually reducing the smoothing weight and allowing the model to focus on more challenging samples during training. FE2H [333] introduces a document selection module that iteratively performs binary classification tasks to select relevant documents by simply adding a prediction layer on a PLM. Then, it trains the reader module on a single-hop QA dataset and transfers it into the multi-hop QA task inspired by humans' progressive learning process.

5.3.3 Open-domain question answering

Task definition: Open-domain question answering methods are required to retrieve relevant passages from a large-scale corpus and generate the final answer based on the retrieved passages. The ODQA task is more challenging than MRC and MHRC, which search the support facts within a smaller set of candidate passages. Recently, ODQA has been widely used to test the retrieval augmented generation (RAG) systems [342, 343].

Most ODQA methods follow a retrieve-and-read pipeline [28, 30, 408]. The objective of the retrieval phase is to retrieve evidence-related passages from a large symbolic knowledge corpus, such as Wikipedia¹⁾. The retriever can be divided into sparse retrieval and dense retrieval. Sparse retrieval methods rely on word-level matching to link vocabulary with documents. Notable methods include Boolean Retrieval [409], BM25 [410], SPLADE [411], and UniCOIL [412]. Dense retrieval methods capture deep semantic information to comprehend the underlying semantics of documents, thereby enhancing retrieval accuracy. Key examples include DPR [413], ANCE [414], RocketQA [415], E5 [416], DrBoost [417], and SimLM [418]. The goal of the reading phase is comprehension and reasoning, akin to MRC, to derive answers based on the retrieved passages. Generally, existing readers can be categorized into extractive readers and generative readers. Extractive readers predict an answer span from the retrieved passages. Notable methods include REALM [49], Skylinebuilderretro [50], RETRO [51], and BPR [52]. Generative Readers generate answers in natural language using sequence-to-sequence models. Key examples include RAG [66], Fusion-in-Decoder [73], MDR [419], and RALM [420].

Recently, the emergence of LLMs has demonstrated their potential for open-domain question answering [336]. This section mainly investigates how to leverage the LLMs to optimize the retrieve-and-read pipeline. At the retrieval stage, the LLMs can be utilized for query reformulation and Enhanced retrieval. At the reader stage, LLMs can serve as augmented readers.

5.3.3.1 Query reformulation methods

Query reformulation methods focus on refining input questions to convey user intent more accurately. For instance, EAR [334] first applies a query expansion model to generate a diverse set of queries and then uses a query reranker to select the ones that could lead to better retrieval results. QPaug [335] decomposes the original questions into multiple-step sub-questions. By augmenting the original question with detailed sub-questions and planning, QPaug can make the query more specific on what needs to be retrieved, improving the retrieval performance. Chain-of-Rewrite [336] finds that current methods face challenges stemming from term mismatch and limited interaction between information retrieval systems and LLMs. Hence, it leverages the guidance and feedback gained from the analysis to provide faithful and consistent extensions for effective question answering. Specifically, CLASS [80] employs LLMs for query transformation via in-context learning in Cross-lingual ODQA tasks.

5.3.3.2 Enhanced retrieval methods

Enhanced retrieval methods adopt LLMs as knowledge sources to provide relevant contextual documents, thereby increasing the likelihood of uncovering the correct answer. This method category is

¹⁾<https://www.wikipedia.org>

special because they directly use LLMs to generate evidence-related passages, replacing the retrieval process. Although they do not directly utilize the knowledge from a symbolic knowledge base, we still introduce it due to its advanced features. For instance, HintQA [337] produces multiple hints for each question. Then, it substitutes the retrieved passages and generated contexts with the generated hints. GenRead [338] prompts an LLM to generate contextual documents based on a given question and then reads the generated documents to produce the final answer. Self-Prompting [339] prompts LLMs step by step to generate multiple pseudo QA pairs with background passages and explanations entirely from scratch. These generated elements are then utilized for in-context learning. MedGENIE [340] prompts a medical LLM to furnish multi-view background contexts for a given question. Then, it designs two readers for prompting LLMs and fine-tuning SLMs, respectively.

5.3.3.3 *Augmented reader methods*

At the reader stage, LLMs can serve as augmented readers, effectively minimizing distractions from irrelevant documents and improving the quality of the context. To overcome the challenges posed by irrelevant retrieved documents and overconfident scores, DAS [341] propose a negation-based instruction to allow LLMs to abstain from answering. Then, it designs a score adjustment strategy to adjust the answer scores by reflecting the query generation score as the relevance between the given query-document pairs. Considering that LLMs cannot precisely assess the relevance of retrieved documents, thus likely leading to misleading or even incorrect utilization of external knowledge, REAR [342] incorporates an assessment module that precisely assesses the relevance of retrieved documents and proposes an improved training method based on bi-granularity relevance fusion and noise-resistant training. RE-RAG [343] introduces a relevance estimator that not only provides relative relevance between contexts as previous rerankers did but also provides confidence, which can be used to classify whether the given context is helpful in answering the given question. FastFiD [344] performs sentence selection post the output of the reader's encoder and maintains only the essential sentences as references for the reader's decoder, thereby significantly reducing the inference time for each query. Considering that the retrieved context may contain noise and irrelevant information and augmenting noisy context can potentially distract LLMs, TA-ARE [345] dynamically determines retrieval necessity and relies only on LLMs' parametric knowledge when deemed unnecessary.

5.3.4 *Multi-document question answering*

Task definition: Multi-document question answering methods aims to find the supporting facts from multiple entire documents. MDQA demands a thorough understanding of the logical associations among the contents and structures of documents.

Although some methods also claim to perform document-based QA, they typically focus on paragraphs with key information, not the entire document. An entire document is usually much longer than a paragraph and contains more distracting information. MDQA requires methods to identify support facts from the entire document, which is challenging. First, an entire document can be very lengthy, and supporting facts may comprise only a tiny part. Moreover, the text within a document is often on a single topic, making different passages highly related and difficult to distinguish.

Recently, a few excellent methods have been proposed for MDQA. KGP [346] formulates the proper context in prompting LLMs for MD-QA, which consists of a graph construction module and a graph traversal module. For graph construction, KGP creates a KG over multiple documents with nodes symbolizing passages or document structures (e.g., pages/tables) and edges denoting the semantic/lexical similarity between passages or document structural relations. For graph traversal, KGP designs an LLM-based graph traversal agent that navigates across nodes and gathers supporting passages to assist LLMs in MD-QA. Considering that some questions often require synthesizing information from multiple frequently unrelated documents, CuriousLLM [347] fine-tunes a decoder-only LLM to emulate the curious nature of a human researcher to generate follow-up questions based on both the initial user query and passages retrieved in previous steps. These questions serve as a guide to identify the most relevant neighboring passages for the subsequent hops in the search process.

5.3.5 Datasets

In this section, we have compiled statistics on some commonly-used datasets related to graph-based reasoning, including (1) Domain: the domain of knowledge corresponding to the datasets; (2) # Ques.: Question number; (3) # Pas.: Passage number; (4) Q&A source: The main construction methods of questions and answers. They are mainly divided into three categories: “Generate”, “Expert”, and “Crowdsourcing”. “Generate” refers to the design of programs for automated generation, “Expert” refers to direct crawling from professional websites or carefully designed by domain experts, and “Crowdsourcing” refers to completion by crowdsourcing workers with general cultural levels; (5) Passage source: The source of passages; (6) Links: The storage address of the datasets. The statistical results are shown in Table 5.

Datasets	Domain	# Ques.	# Pas.	Q&A source	Passages source	Links
CNN/DailyMail [383]	News	1.38 M	312 K	Generate	CNN and DailyMail websites	https://
NewsQA [421]	News	120K	12.7 K	Crowdsourcing	CNN websites	https://
PeopleDaily/CFT [421]	News	880 K	60K	Generate	People Daily websites	https://
TriviaQA [385]	News	95.9K	663 K	Crowdsourcing	Bing Search	https://
RACE [386]	Science	97.6 K	28 K	Expert	English Exam	https://
ARC [422]	Science	7,787	14 M	Expert	Science Exam	https://
SQuAD1.1 [384]	General	107.8 K	536	Crowdsourcing	Wikipedia	https://
SQuAD2.0 [423]	General	150 K	536	Crowdsourcing	Wikipedia	https://
WikiQA [424]	General	3,047	29.3 K	Crowdsourcing	Wikipedia	https://
SearchQA [425]	General	140 K	6.9 M	Crowdsourcing	Google Search	https://
HotpotQA [426]	General	113 K	-	Crowdsourcing	Wikipedia	https://
Natural Questions [427]	General	3.09 M	323 K	Crowdsourcing	Google Search	https://
2WikiMultiHopQA [428]	General	192.6 K	-	Crowdsourcing	Wikipedia and Wikidata	https://
IIRC [429]	General	13 K	-	Crowdsourcing	Wikipedia	https://
FanOutQA [430]	General	8,339	-	Crowdsourcing	Wikipedia	https://

Table 5 Dataset statistics of text-based reasoning.

5.4 Heterogeneous reasoning

Graph-based reasoning, table-based reasoning, and text-based reasoning have all been individually studied extensively. However, reasoning based on two or more heterogeneous symbolic knowledge bases, known as heterogeneous question answering (Heterogeneous QA), is under-studied [353]. Exploring how to explore the knowledge from multiple heterogeneous symbolic knowledge bases fully is extremely important for enhancing the practicality of reasoning methods.

5.4.1 Heterogeneous question answering

Task definition: Heterogeneous question answering methods aim to find the evidence from heterogeneous knowledge bases to answer a knowledge-intensive question.

Some methods investigate how to leverage symbolic knowledge from different sources, including those on closed domain [302, 431–434] and open domain [435–438], but very limited existing work experiments on graph, table, and text, simultaneously. To promote relevant research, CONVMIX [350] and COMPMIX [439] collect the heterogeneous QA datasets that require knowledge from all three heterogeneous sources. A simple solution for handling heterogeneous QA is to assemble several specialized systems. In this approach, the input question is dispatched to multiple sub-systems, and one of them is chosen to provide the final answer. Although this method can leverage state-of-the-art models optimized for various information sources, it significantly increases the complexity of the entire system. Additionally, it poses challenges in addressing questions that require reasoning across multiple sources of information [348]. Hence, constructing an integrated system compatible with multiple heterogeneous symbolic knowledge bases is essential and promising. Current methods with integrated systems can be divided into structured-unified and human-imitated methods.

5.4.1.1 Structure-unified methods

The structured-unified methods convert multiple heterogeneous bases into one type. The first type of structured-unified method converts different structures to unstructured text, and the second one converts

different structures to structured graphs. In the first type of method, UniK-QA [348] flatten the lists, tables, and KGs to text using simple heuristics methods. Then, it adopts a text-based QA method as the solution to make full use of the powerful PLMs. UDT-QA [349] unifies both representation and model for ODQA over structured data and unstructured text. The key idea is to augment the retriever with a data-to-text verbalizer for accessing heterogeneous knowledge bases, i.e., KGs from WikiData, tables and texts from Wikipedia. Convinse [350] learns an explicit and structured representation of an incoming question and its conversational context. It harnesses this frame-like representation to uniformly capture relevant evidence from KB, text, and tables. Finally, it adopts a fusion-in-decoder model to generate the answer.

In the second type of method, TrustUQA [351] designs a condition graph to unify tables and KGs and uses an LLM and demonstration-based two-level method for reasoning on condition graph. Explaignn [352] constructs a heterogeneous graph from entities and evidence snippets retrieved from a KG, a text corpus, web tables, and infoboxes. This large graph is then iteratively reduced via GNNs incorporating question-level attention until the best answers and explanations are distilled out. The former gives up the advantage of using formal query languages on structured data, which can support operations such as ranking and averaging. The latter gives up the advantage of the expressiveness and versatility of free-text knowledge representation. As [353] points out, the first type of method sacrifices the benefits of using formal query languages on structured data, which can support operations like ranking and averaging. The second type of method relinquishes the expressiveness and versatility offered by free-text knowledge representation.

5.4.1.2 Human-imitated methods

The human-imitated methods integrate reasoning steps over heterogeneous knowledge bases by mimicking how humans find responses to questions, which break down QA solution processes as tool calls and thoughts. For instance, HumanIQ [81] proposes a human-like approach that teaches LLMs to gather heterogeneous information by imitating how humans use retrieval tools. During the preparation stage, the method is required to identify suitable tools and solution processes using those tools. Then, it leverages an LLM to replicate the solution processes at the inference stage. Similarly, SPAGHETTI [353] obtain evidence from heterogeneous sources in parallel, including structured KG, plain text, linearized tables, infoboxes, and LLM-generated claims that are verified, and gather those evidence to generate the final answer using a few-shot LLM.

5.4.2 Datasets

In this section, we have compiled statistics on some commonly-used datasets related to heterogeneous reasoning, including (1) Domain: the domain of knowledge corresponding to the datasets; (2) # Ques.: Question number; (3) Q&A source: The main construction methods of questions and answers. They are mainly divided into three categories: “Generate”, “Expert”, and “Crowdsourcing”. “Generate” refers to the design of programs for automated generation, “Expert” refers to direct crawling from professional websites or carefully designed by domain experts, and “Crowdsourcing” refers to completion by crowdsourcing workers with general cultural levels; (4) KG: Whether they use KG as knowledge bases; (5) KG size: the number of triplets in KG used; (6) Text: Whether they use Text as knowledge bases; (7) # Pas.: the number of passages; (8) Table: Whether they use Table as knowledge bases; (8) # Table: The number of tables; (9) OR: Whether they support open retrieval; (10) HQ: Whether the questions are constructed by human; (11) OD: Whether the answers could be found in open domain; (12) Links: The storage address of the datasets. The statistical results are shown in Table 6.

6 Future directions

To enhance the feasibility of reasoning systems in real-world applications, we may need to focus on three key aspects: performance, cost, and security. We have provided some insights into future directions based on these aspects to foster the development of reasoning systems.

Datasets	Domain	# Ques.	Q&A source	KG	KG size	Text	# Pas.	Table	# Table	OR	HQ	OD	Links
WIKIMOVIES [431]	Movie	100 K	Generate	✓	-	✓	17 K	✗	-	✓	✗	✗	https://
HYBRIDQA [432]	General	70 K	Crowdsourcing	✗	-	✓	293 K	✓	13 K	✗	✓	✓	https://
MULTIMODALQA [440]	General	30 K	Generate	✗	-	✓	218 K	✓	10 K	✓	✗	✓	https://
OTT-QA [435]	General	45 K	Crowdsourcing	✗	-	✓	-	✓	-	✓	✓	✓	https://
MANYMODALQA [441]	General	10 K	Crowdsourcing	✗	-	✓	3,789	✓	528	✗	✓	✓	https://
TAT-QA [442]	General	17 K	Crowdsourcing	✗	-	✓	3,902	✓	7,431	✗	✓	✗	https://
FINQA [443]	Finance	8,281	Crowdsourcing	✗	-	✓	-	✓	-	✗	✓	✗	https://
HETPQA [444]	Product	6,000	Crowdsourcing	✗	-	✓	-	✓	-	✗	✓	✗	https://
COMPPIX [439]	General	9,410	Crowdsourcing	✓	-	✓	-	✓	-	✓	✓	✓	https://

Table 6 Dataset statistics of heterogeneous reasoning.

6.1 Improving OOD generalization and robustness of reasoning

Despite the good performance and generalization of LLMs in various general reasoning tasks, such as mathematical reasoning and commonsense reasoning, in some specific domains like healthcare and biology, we still rely on domain-specific data to construct the required reasoning systems based on LLMs. This includes introducing domain-specific pre-training corpora for continued pre-training [445], introducing post-training data for instruction fine-tuning or alignment fine-tuning [446], or incorporating domain-specific knowledge bases for retrieval-augmented reasoning. This poses a hidden risk: due to the limited availability of domain-specific data, these reasoning systems may perform well on seen data but demonstrate poor generalization and instability on unseen out-of-distribution (OOD) instances [447–449]. Therefore, exploring and enhancing the generalization and robustness of these domain-specific reasoning systems is an important research direction. The following lists some potential research directions:

- **Scalable reasoning training data construction:** Scalable data construction methods based on limited in-domain data may help reasoning systems generalize to more reasoning examples since it can expand the distribution encountered during the training process [450, 451]. In this direction, ensuring the quality and diversity of the expanded data will be a core research issue.
- **Rule-guided reasoning:** One of the advantages of symbolic reasoning [452] is that it is not limited by the finite distribution of training data, which may help improve the generalization and robustness of reasoning systems. Specifically, future research could explore how to derive executable reasoning rules from existing in-domain data and how to apply these reasoning rules during the inference phase for symbolic reasoning, thereby enhancing its reliability.

6.2 Cost-efficient reasoning

Existing strategies to enhance reasoning system performance, such as chain-of-thought, self-consistency, and feedback-enhanced reasoning, often become inefficient due to the need to generate more tokens. In practical applications, we desire reasoning systems that not only demonstrate good performance but also have minimal reasoning time and cost. Hence, future research could focus on proposing cost-efficient reasoning strategies while improving or maintaining reasoning system performance. To address this issue, future research can attempt the following:

- **Incorporating symbolic reasoning strategies:** Although rule-based symbolic systems [452, 453] and probabilistic statistical models are less effective compared to complex neural systems like LLMs, they can handle simple reasoning problems or some simple reasoning steps in complex reasoning problems. Exploring how to integrate symbolic systems, statistical models, and LLMs to collaboratively reason may help achieve cost-efficient reasoning while maintaining powerful performance.
- **Efficiency-adaptive reasoning:** For humans, the reasoning cost required for simple problems differs from that for more complex problems. How to adaptively engage in fast-thinking reasoning for simpler problems and slow-thinking reasoning for more difficult ones is a promising approach to achieving cost-efficient reasoning [454, 455].

6.3 Ensuring the safety of reasoning

Ensuring the safety of reasoning involves two aspects. First, ensuring that the reasoning processes and results do not contain harmful or potentially dangerous contents. Second, ensuring that private information from local symbolic knowledge bases is not leaked to cloud-based LLMs.

6.3.1 Preventing toxicity while reasoning

Reasoning is a knowledge-intensive task where we infer new knowledge from existing knowledge. We want this new knowledge to be harmless and not cause negative impacts on society or individuals. However, current toxicity detection or detoxification work mainly focuses on hate or biased speech in general texts [456–460], not on such knowledge-intensive reasoning tasks. Compared to general texts, texts in reasoning processes and results contain more specialized knowledge, such as chemistry, biology, physics, and mathematics. For example, designing molecular formulas for drug creation can be harmful but hard to detect based on text semantics alone. Therefore, conducting toxicity detection or detoxification for such knowledge-intensive reasoning tasks is a challenge and a significant concern. The following is a potential research direction:

- **Incorporating domain-specific expertise:** To ensure the security of reasoning processes and results, augmenting traditional toxicity detection approaches with domain-specific expertise may be necessary. This involves integrating interdisciplinary insights from fields like chemistry, biology, and physics with AI safety frameworks to identify potential harm beyond surface-level text semantics.

6.3.2 Protecting data privacy while reasoning

Recently, an increasing number of applications have integrated their local symbolic knowledge bases with third-party LLMs for retrieval-augmented generation [461], which enhances the accuracy and credibility of the generation for domain-specific reasoning. However, it poses potential privacy leakage risks. Taking KGQA as an example, the local KG contains many triplets involving private data, especially in entities. Meanwhile, many high-performance LLMs in the cloud operate in a black-box manner through API calls. Current methods, particularly IR-based methods, require these triplets to be transmitted indiscriminately to the LLM, even though most entities are unnecessary for the reasoning process. Therefore, protecting data privacy while reasoning in a de-identification scenario by combining the local knowledge base with third-party LLMs is a challenge and an important focus. Here is a potential research direction:

- **Reasoning in scenarios where most entity names are anonymized:** To ensure that private information in the local knowledge base is not exposed to third-party LLM during the reasoning process, it may be necessary to anonymize the specific names of entities in the local knowledge base, such as by replacing them with entity IDs. During the reasoning process, LLM would have to rely solely on the types of entities and their relations, while the semantics of the entities and the rich background knowledge stored in the LLM would no longer be available. This involves the application of traditional privacy computation methods and fully utilizing the LLM’s logical reasoning capabilities.

7 Conclusion

In this paper, we provide a comprehensive survey on reasoning methods with a specific focus on the usage of knowledge bases, addressing a gap in existing literature. By categorizing knowledge bases into symbolic and parametric types, we offer a novel perspective on how reasoning can be enhanced when leveraging different formats of stored information. Symbolic knowledge bases, such as KGs and tables, offer explicit and human-readable knowledge, while parametric knowledge bases encode knowledge implicitly within parameters, such as large language models. Then, we investigate how these knowledge bases, individually and in combination, support reasoning processes. Additionally, this survey proposes some potential future directions in reasoning research. By providing this comprehensive overview, we hope to inspire further exploration and advancements in the field, ultimately contributing to developing artificial intelligence systems with more robust reasoning capabilities.

References

- 1 Fei Yu, Hongbo Zhang, Prayag Tiwari, and Benyou Wang. Natural language reasoning, a survey. *ACM Computing Surveys*, 56(12):1–39, 2024.
- 2 Jie Huang and Kevin Chen-Chuan Chang. Towards reasoning in large language models: A survey. In *Findings of the Association for Computational Linguistics: ACL 2023*, pages 1049–1065, 2023.
- 3 Maggi Banning. Clinical reasoning and its application to nursing: concepts and research studies. *Nurse education in practice*, 8(3):177–183, 2008.
- 4 Erwin B Montgomery Jr. *Medical reasoning: the nature and use of medical knowledge*. Oxford University Press, 2018.
- 5 Travis Zack, Gurpreet Dhaliwal, Rabih Geha, Mary Margaretten, Sara Murray, and Julian C Hong. A clinical reasoning-encoded case library developed through natural language processing. *Journal of General Internal Medicine*, 38(1):5–11, 2023.
- 6 Changyu Chen, Xiting Wang, Ting-En Lin, Ang Lv, Yuchuan Wu, Xin Gao, Ji-Rong Wen, Rui Yan, and Yongbin Li. Masked thought: Simply masking partial reasoning steps can improve mathematical reasoning learning of language models. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 5872–5900. Association for Computational Linguistics, 2024.
- 7 Aitor Lewkowycz, Anders Andreassen, David Dohan, Ethan Dyer, Henryk Michalewski, Vinay Ramasesh, Ambrose Slone, Cem Anil, Imanol Schlag, Theo Gutman-Solo, et al. Solving quantitative reasoning problems with language models. *Advances in Neural Information Processing Systems*, 35:3843–3857, 2022.
- 8 Wenhu Chen, Xueguang Ma, Xinyi Wang, and William W Cohen. Program of thoughts prompting: Disentangling computation from reasoning for numerical reasoning tasks. *Transactions on Machine Learning Research*.
- 9 Guijin Son, Hanearl Jung, Moonjeong Hahm, Keonju Na, and Sol Jin. Beyond classification: Financial reasoning in state-of-the-art language models. In *Proceedings of the Fifth Workshop on Financial Technology and Natural Language Processing and the Second Multimodal AI For Financial Forecasting*, pages 34–44, 2023.
- 10 Ziqiang Yuan, Kaiyuan Wang, Shoutai Zhu, Ye Yuan, Jingya Zhou, Yanlin Zhu, and Wenqi Wei. Finllms: A framework for financial reasoning dataset generation with large language models. *arXiv preprint arXiv:2401.10744*, 2024.
- 11 Yilun Zhao, Yitao Long, Hongjun Liu, Ryo Kamoi, Linyong Nan, Lyuhao Chen, Yixin Liu, Xiangru Tang, Rui Zhang, and Arman Cohan. Docmath-eval: Evaluating math reasoning capabilities of llms in understanding financial documents. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 16103–16120, 2024.
- 12 Nan Duan, Duyu Tang, and Ming Zhou. Machine reasoning: Technology, dilemma and future. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing: Tutorial Abstracts*, pages 1–6, 2020.
- 13 Siyuan Wang, Zhongkun Liu, Wanjun Zhong, Ming Zhou, Zhongyu Wei, Zhumin Chen, and Nan Duan. From lsat: The progress and challenges of complex reasoning. *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, 30:2201–2216, 2022.
- 14 Prajjwal Bhargava and Vincent Ng. Commonsense knowledge reasoning and generation with pre-trained language models: A survey. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 36, pages 12317–12325, 2022.
- 15 Shuofei Qiao, Yixin Ou, Ningyu Zhang, Xiang Chen, Yunzhi Yao, Shumin Deng, Chuanqi Tan, Fei Huang, and Huajun Chen. Reasoning with language model prompting: A survey. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 5368–5393, 2023.
- 16 Daniel C Krawczyk. The cognition and neuroscience of relational reasoning. *Brain research*, 1428:13–23, 2012.
- 17 Peter C Wason. Reasoning about a rule. *Quarterly journal of experimental psychology*, 20(3):273–281, 1968.
- 18 Nicholas Rescher. *Philosophical reasoning: A study in the methodology of philosophizing*. 2001.
- 19 John Arthur Passmore. *Philosophical reasoning*. 1961.
- 20 M Huth. *Logic in Computer Science: Modelling and reasoning about systems*. Cambridge University Press, 2004.
- 21 Aske Plaat, Annie Wong, Suzan Verberne, Joost Broekens, Niki van Stein, and Thomas Back. Reasoning with large language models, a survey. *arXiv preprint arXiv:2407.11511*, 2024.
- 22 Ke Liang, Lingyuan Meng, Meng Liu, Yue Liu, Wenxuan Tu, Siwei Wang, Sihang Zhou, Xinwang Liu, Fuchun Sun, and Kunlun He. A survey of knowledge graph reasoning on graph types: Static, dynamic, and multi-modal. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2024.
- 23 Xiaojun Chen, Shengbin Jia, and Yang Xiang. A review: Knowledge reasoning over knowledge graph. *Expert systems with applications*, 141:112948, 2020.
- 24 Siyuan Meng, Jie Zhou, Xuxin Chen, Yufei Liu, Fengyuan Lu, and Xinli Huang. Structure-information-based reasoning over the knowledge graph: A survey of methods and applications. *ACM Transactions on Knowledge Discovery from Data*, 18(8):1–42, 2024.
- 25 Lauren Nicole DeLong, Ramon Fernández Mir, and Jacques D Fleuriot. Neurosymbolic ai for reasoning over knowledge graphs: A survey. *IEEE Transactions on Neural Networks and Learning Systems*, 2024.
- 26 Yunshi Lan, Gaole He, Jinhao Jiang, Jing Jiang, Wayne Xin Zhao, and Ji-Rong Wen. A survey on complex knowledge base question answering: Methods, challenges and solutions. In *Proceedings of the Thirtieth International Joint Conference on Artificial Intelligence*, 2021.
- 27 Miao Su, ZiXuan Li, Zhuo Chen, Long Bai, Xiaolong Jin, and Jiafeng Guo. Temporal knowledge graph question answering: A survey. *arXiv preprint arXiv:2406.14191*, 2024.
- 28 Fengbin Zhu, Wenqiang Lei, Chao Wang, Jianming Zheng, Soujanya Poria, and Tat-Seng Chua. Retrieving and reading: A comprehensive survey on open-domain question answering. *arXiv preprint arXiv:2101.00774*, 2021.
- 29 Xiaoyu Shen, Svitlana Vakulenko, Marco Del Tredici, Gianni Barlacchi, Bill Byrne, and Adrià de Gispert. Low-resource dense retrieval for open-domain question answering: A comprehensive survey. *arXiv preprint arXiv:2208.03197*, 2022.
- 30 Qin Zhang, Shangsi Chen, Dongkuan Xu, Qingqing Cao, Xiaojun Chen, Trevor Cohn, and Meng Fang. A survey for efficient open domain question answering. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 14447–14465, 2023.
- 31 Ronald Fagin, Joseph Y Halpern, Yoram Moses, and Moshe Vardi. *Reasoning about knowledge*. MIT press, 2004.
- 32 Herbert A Simon and Allen Newell. Computer science as empirical inquiry: symbols and search. *Communications of the ACM*, 19(3):11–126, 1976.
- 33 David E Rumelhart, Geoffrey E Hinton, and Ronald J Williams. Learning representations by back-propagating errors. *nature*, 323(6088):533–536, 1986.
- 34 Randall Davis, Howard Shrobe, and Peter Szolovits. What is a knowledge representation? *AI magazine*, 14(1):17–17, 1993.
- 35 Yann LeCun, Yoshua Bengio, and Geoffrey Hinton. Deep learning. *nature*, 521(7553):436–444, 2015.
- 36 Ke Liang, Lingyuan Meng, Meng Liu, Yue Liu, Wenxuan Tu, Siwei Wang, Sihang Zhou, Xinwang Liu, Fuchun Sun, and Kunlun He. A survey of knowledge graph reasoning on graph types: Static, dynamic, and multi-modal. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 46(12):9456–9478, 2024.
- 37 Christine Vanoirbeek. Formatting structured tables. In *EP92 (Proceedings of Electronic Publishing, 1992)*, pages 291–309. Cambridge University Press UK, 1992.
- 38 Yuan Sui, Mengyu Zhou, Mingjie Zhou, Shi Han, and Dongmei Zhang. Table meets llm: Can large language models

- understand structured table data? a benchmark and empirical study. In *Proceedings of the 17th ACM International Conference on Web Search and Data Mining*, pages 645–654, 2024.
- 39 John Dagdelen, Alexander Dunn, Sanghoon Lee, Nicholas Walker, Andrew S Rosen, Gerbrand Ceder, Kristin A Persson, and Anubhav Jain. Structured information extraction from scientific text with large language models. *Nature Communications*, 15(1):1418, 2024.
- 40 Alon Talmor, Oyvind Tafjord, Peter E. Clark, Yoav Goldberg, and Jonathan Berant. Leap-of-thought: Teaching pre-trained models to systematically reason over implicit knowledge. *Neural Information Processing Systems, Neural Information Processing Systems*, Jan 2020.
- 41 Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. BERT: Pre-training of deep bidirectional transformers for language understanding. In Jill Burstein, Christy Doran, and Thamar Solorio, editors, *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 4171–4186, Minneapolis, Minnesota, June 2019. Association for Computational Linguistics.
- 42 Yinhan Liu. Roberta: A robustly optimized bert pretraining approach. *arXiv preprint arXiv:1907.11692*, 364, 2019.
- 43 Zhenzhong Lan, Mingda Chen, Sebastian Goodman, Kevin Gimpel, Piyush Sharma, and Radu Soricut. Albert: A lite bert for self-supervised learning of language representations. In *International Conference on Learning Representations*, 2020.
- 44 Pengcheng He, Xiaodong Liu, Jianfeng Gao, and Weizhu Chen. {DEBERTA}: {DECODING}-{enhanced} {bert} {with} {disentangled} {attention}. In *International Conference on Learning Representations*, 2021.
- 45 Alexis Conneau and Guillaume Lample. Cross-lingual language model pretraining. *Advances in neural information processing systems*, 32, 2019.
- 46 Zhilin Yang, Zihang Dai, Yiming Yang, Jaime Carbonell, Ruslan Salakhutdinov, and Quoc V. Le. *XLNet: generalized autoregressive pretraining for language understanding*. Curran Associates Inc., Red Hook, NY, USA, 2019.
- 47 Li Dong, Nan Yang, Wenhui Wang, Furu Wei, Xiaodong Liu, Yu Wang, Jianfeng Gao, Ming Zhou, and Hsiao-Wuen Hon. Unified language model pre-training for natural language understanding and generation. *Advances in neural information processing systems*, 32, 2019.
- 48 Kevin Clark, Minh-Thang Luong, Quoc V. Le, and Christopher D. Manning. ELECTRA: pre-training text encoders as discriminators rather than generators. In *8th International Conference on Learning Representations, ICLR 2020, Addis Ababa, Ethiopia, April 26-30, 2020*. OpenReview.net, 2020.
- 49 Kelvin Guu, Kenton Lee, Zora Tung, Panupong Pasupat, and Mingwei Chang. Retrieval augmented language model pre-training. In *International conference on machine learning*, pages 3929–3938. PMLR, 2020.
- 50 Yuxiang Wu, Pasquale Minervini, Pontus Stenetorp, and Sebastian Riedel. Don't read too much into it: Adaptive computation for open-domain question answering. In *Proceedings of SustaiNLP: Workshop on Simple and Efficient Natural Language Processing*, pages 63–72, 2020.
- 51 Sebastian Borgeaud, Arthur Mensch, Jordan Hoffmann, Trevor Cai, Eliza Rutherford, Katie Millican, George Bm Van Den Driessche, Jean-Baptiste Lespiau, Bogdan Damoc, Aidan Clark, et al. Improving language models by retrieving from trillions of tokens. In *International conference on machine learning*, pages 2206–2240. PMLR, 2022.
- 52 Ikuya Yamada, Akari Asai, and Hannaneh Hajishirzi. Efficient passage retrieval with hashing for open-domain question answering. In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 2: Short Papers)*, pages 979–986. Association for Computational Linguistics.
- 53 Zhuosheng Zhang, Junjie Yang, and Hai Zhao. Retrospective reader for machine reading comprehension. In *Proceedings of the AAAI conference on artificial intelligence*, volume 35, pages 14506–14514, 2021.
- 54 Chen Zhang, Yuxuan Lai, Yansong Feng, and Dongyan Zhao. Extract, integrate, compete: Towards verification style reading comprehension. In *Findings of the Association for Computational Linguistics: EMNLP 2021*, pages 2976–2986, 2021.
- 55 Zixian Huang, Jiaying Zhou, Chenxu Niu, and Gong Cheng. Spans, not tokens: A span-centric model for multi-span reading comprehension. In *Proceedings of the 32nd ACM International Conference on Information and Knowledge Management*, pages 874–884, 2023.
- 56 Siru Ouyang, Zhuosheng Zhang, and Hai Zhao. Fact-driven logical reasoning for machine reading comprehension. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 38, pages 18851–18859, 2024.
- 57 Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*, 2023.
- 58 Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shrutit Bhosale, et al. Llama 2: Open foundation and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*, 2023.
- 59 Jinze Bai, Shuai Bai, Yunfei Chu, Zeyu Cui, Kai Dang, Xiaodong Deng, Yang Fan, Wenbin Ge, Yu Han, Fei Huang, et al. Qwen technical report. *arXiv preprint arXiv:2309.16609*, 2023.
- 60 Albert Q Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier, et al. Mistral 7b. *arXiv preprint arXiv:2310.06825*, 2023.
- 61 Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Brian Ichter, Fei Xia, Ed H. Chi, Quoc V. Le, and Denny Zhou. Chain-of-thought prompting elicits reasoning in large language models. In *Proceedings of the 36th International Conference on Neural Information Processing Systems, NIPS '22*, Red Hook, NY, USA, 2024. Curran Associates Inc.
- 62 Aman Madaan, Niket Tandon, Prakhar Gupta, Skyler Hallinan, Luyu Gao, Sarah Wiegrefe, Uri Alon, Nouha Dziri, Shrimai Prabhumoye, Yiming Yang, Shashank Gupta, Bodhisatwa Prasad Majumder, Katherine Hermann, Sean Welleck, Amir Yazdanbakhsh, and Peter Clark. Self-refine: Iterative refinement with self-feedback. In Alice Oh, Tristan Naumann, Amir Globerson, Kate Saenko, Moritz Hardt, and Sergey Levine, editors, *Advances in Neural Information Processing Systems 36: Annual Conference on Neural Information Processing Systems 2023, NeurIPS 2023, New Orleans, LA, USA, December 10 - 16, 2023*, 2023.
- 63 Maxwell Nye, Anders Johan Andreassen, Guy Gur-Ari, Henryk Michalewski, Jacob Austin, David Bieber, David Dohan, Aitor Lewkowycz, Maarten Bosma, David Luan, Charles Sutton, and Augustus Odena. Show your work: Scratchpads for intermediate computation with language models. In *Deep Learning for Code Workshop*, 2022.
- 64 Ranjita Naik, Varun Chandrasekaran, Mert Yuksekgonul, Hamid Palangi, and Besmira Nushi. Diversity of thought improves reasoning abilities of llms, 2024.
- 65 Yilun Du, Shuang Li, Antonio Torralba, Joshua B. Tenenbaum, and Igor Mordatch. Improving factuality and reasoning in language models through multiagent debate. In *Forty-first International Conference on Machine Learning*, 2024.
- 66 Patrick Lewis, Ethan Perez, Aleksandra Piktus, Fabio Petroni, Vladimir Karpukhin, Naman Goyal, Heinrich Küttler, Mike Lewis, Wen-tau Yih, Tim Rocktäschel, et al. Retrieval-augmented generation for knowledge-intensive nlp tasks. *Advances in Neural Information Processing Systems*, 33:9459–9474, 2020.
- 67 Dan Hendrycks, Collin Burns, Saurav Kadavath, Akul Arora, Steven Basart, Eric Tang, Dawn Song, and Jacob Steinhardt. Measuring mathematical problem solving with the math dataset. In *Thirty-fifth Conference on Neural Information Processing Systems Datasets and Benchmarks Track (Round 2)*, 2021.

- 68 Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J Liu. Exploring the limits of transfer learning with a unified text-to-text transformer. *Journal of machine learning research*, 21(140):1–67, 2020.
- 69 M Lewis. Bart: Denoising sequence-to-sequence pre-training for natural language generation, translation, and comprehension. *arXiv preprint arXiv:1910.13461*, 2019.
- 70 Y Liu. Multilingual denoising pre-training for neural machine translation. *arXiv preprint arXiv:2001.08210*, 2020.
- 71 Kaitao Song. Mass: Masked sequence to sequence pre-training for language generation. *arXiv preprint arXiv:1905.02450*, 2019.
- 72 Team GLM, Aohan Zeng, Bin Xu, Bowen Wang, Chenhui Zhang, Da Yin, Diego Rojas, Guanyu Feng, Hanlin Zhao, Hanyu Lai, et al. Chatglm: A family of large language models from glm-130b to glm-4 all tools. *arXiv preprint arXiv:2406.12793*, 2024.
- 73 Gautier Izacard and Édouard Grave. Leveraging passage retrieval with generative models for open domain question answering. In *Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics: Main Volume*, pages 874–880, 2021.
- 74 Antoine Bordes, Nicolas Usunier, Alberto Garcia-Durán, Jason Weston, and Oksana Yakhnenko. Translating embeddings for modeling multi-relational data. In *Proceedings of the 26th International Conference on Neural Information Processing Systems - Volume 2*, NIPS’13, page 2787–2795, Red Hook, NY, USA, 2013. Curran Associates Inc.
- 75 Tingsong Jiang, Tianyu Liu, Tao Ge, Lei Sha, Sujian Li, Baobao Chang, and Zhifang Sui. Encoding temporal information for time-aware link prediction. In Jian Su, Kevin Duh, and Xavier Carreras, editors, *Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing*, pages 2350–2354, Austin, Texas, November 2016. Association for Computational Linguistics.
- 76 Xuezhong Wang, Jason Wei, Dale Schuurmans, Quoc V. Le, Ed H. Chi, Sharan Narang, Aakanksha Chowdhery, and Denny Zhou. Self-consistency improves chain of thought reasoning in language models. In *The Eleventh International Conference on Learning Representations, ICLR 2023, Kigali, Rwanda, May 1-5, 2023*. OpenReview.net, 2023.
- 77 Tian Liang, Zhiwei He, Wenxiang Jiao, Xing Wang, Yan Wang, Rui Wang, Yujie Yang, Shuming Shi, and Zhaopeng Tu. Encouraging divergent thinking in large language models through multi-agent debate. In Yaser Al-Onaizan, Mohit Bansal, and Yun-Nung Chen, editors, *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, pages 17889–17904, Miami, Florida, USA, November 2024. Association for Computational Linguistics.
- 78 Shulin Cao, Jiaxin Shi, Liangming Pan, Lunyu Nie, Yutong Xiang, Lei Hou, Juanzi Li, Bin He, and Hanwang Zhang. KQA pro: A dataset with explicit compositional programs for complex question answering over knowledge base. In Smaranda Muresan, Preslav Nakov, and Aline Villavicencio, editors, *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 6101–6119, Dublin, Ireland, May 2022. Association for Computational Linguistics.
- 79 Mingyu Zheng, Yang Hao, Wenbin Jiang, Zheng Lin, Yajuan Lyu, QiaoQiao She, and Weiping Wang. IM-TQA: A Chinese table question answering dataset with implicit and multi-type table structures. In Anna Rogers, Jordan Boyd-Graber, and Naoaki Okazaki, editors, *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 5074–5094, Toronto, Canada, July 2023. Association for Computational Linguistics.
- 80 Fan Jiang, Tom Drummond, and Trevor Cohn. Pre-training cross-lingual open domain question answering with large-scale synthetic supervision. In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, pages 13906–13933. Association for Computational Linguistics, 2024.
- 81 Jens Lehmann, Dhananjay Bhandiwad, Preetam Gattogi, and Sahar Vahdati. Beyond boundaries: A human-like approach for question answering over structured and unstructured information sources. *Transactions of the Association for Computational Linguistics*, 12:786–802, 2024.
- 82 Zhen Wang, Jianwen Zhang, Jianlin Feng, and Zheng Chen. Knowledge graph embedding by translating on hyperplanes. In *Proceedings of the Twenty-Eighth AAAI Conference on Artificial Intelligence*, AAAI’14, page 1112–1119. AAAI Press, 2014.
- 83 Yankai Lin, Zhiyuan Liu, Maosong Sun, Yang Liu, and Xuan Zhu. Learning entity and relation embeddings for knowledge graph completion. *Proceedings of the AAAI Conference on Artificial Intelligence*, 29(11), February 2015.
- 84 Guoliang Ji, Shizhu He, Liheng Xu, Kang Liu, and Jun Zhao. Knowledge graph embedding via dynamic mapping matrix. In Chengqing Zong and Michael Strube, editors, *Proceedings of the 53rd Annual Meeting of the Association for Computational Linguistics and the 7th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pages 687–696, Beijing, China, July 2015. Association for Computational Linguistics.
- 85 Zhiqing Sun, Zhi-Hong Deng, Jian-Yun Nie, and Jian Tang. Rotate: Knowledge graph embedding by relational rotation in complex space. In *International Conference on Learning Representations*, 2019.
- 86 Ivana Balazević, Carl Allen, and Timothy Hospedales. *Multi-relational poincaré graph embeddings*. Curran Associates Inc., Red Hook, NY, USA, 2019.
- 87 Zhanqiu Zhang, Jianyu Cai, Yongdong Zhang, and Jie Wang. Learning hierarchy-aware knowledge graph embeddings for link prediction. *Proceedings of the AAAI Conference on Artificial Intelligence*, 34(0303):3065–3072, April 2020.
- 88 Maximilian Nickel, Volker Tresp, and Hans-Peter Kriegel. A three-way model for collective learning on multi-relational data. In *Proceedings of the 28th International Conference on International Conference on Machine Learning*, ICML’11, page 809–816, Madison, WI, USA, 2011. Omnipress.
- 89 Bishan Yang, Wen tau Yih, Xiaodong He, Jianfeng Gao, and Li Deng. Embedding entities and relations for learning and inference in knowledge bases. In *International Conference on Learning Representations*, 2014.
- 90 Seyed Mehran Kazemi and David Poole. Simple embedding for link prediction in knowledge graphs. In *Proceedings of the 32nd International Conference on Neural Information Processing Systems*, NIPS’18, page 4289–4300, Red Hook, NY, USA, 2018. Curran Associates Inc.
- 91 Théo Trouillon, Johannes Welbl, Sebastian Riedel, Éric Gaussier, and Guillaume Bouchard. Complex embeddings for simple link prediction. In *Proceedings of the 33rd International Conference on International Conference on Machine Learning - Volume 48*, ICML’16, page 2071–2080. JMLR.org, 2016.
- 92 Shuai Zhang, Yi Tay, Lina Yao, and Qi Liu. *Quaternion knowledge graph embeddings*. Curran Associates Inc., Red Hook, NY, USA, 2019.
- 93 Ivana Balazević, Carl Allen, and Timothy Hospedales. TuckER: Tensor factorization for knowledge graph completion. In Kentaro Inui, Jing Jiang, Vincent Ng, and Xiaojun Wan, editors, *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 5185–5194, Hong Kong, China, November 2019. Association for Computational Linguistics.
- 94 Tim Dettmers, Pasquale Minervini, Pontus Stenetorp, and Sebastian Riedel. Convolutional 2d knowledge graph embeddings. *Proceedings of the AAAI Conference on Artificial Intelligence*, 32(1), April 2018.
- 95 Michael Schlichtkrull, Thomas N. Kipf, Peter Bloem, Rianne van den Berg, Ivan Titov, and Max Welling. Modeling relational data with graph convolutional networks. page 593–607, Berlin, Heidelberg, 2018. Springer-Verlag.
- 96 Luis Antonio Galárraga, Christina Teflioudi, Katja Hose, and Fabian Suchanek. Amie: association rule mining under incomplete evidence in ontological knowledge bases. In *Proceedings of the 22nd International Conference on World Wide*

- Web, WWW '13, page 413–422, New York, NY, USA, 2013. Association for Computing Machinery.
- 97 Christian Meilicke, Manuel Fink, Yanjie Wang, Daniel Ruffinelli, Rainer Gemulla, and Heiner Stuckenschmidt. Fine-grained evaluation of rule- and embedding-based systems for knowledge graph completion. In *The Semantic Web – ISWC 2018: 17th International Semantic Web Conference, Monterey, CA, USA, October 8–12, 2018, Proceedings, Part I*, page 3–20, Berlin, Heidelberg, 2018. Springer-Verlag.
 - 98 Christian Meilicke, Melisachew Wudage Chekol, Daniel Ruffinelli, and Heiner Stuckenschmidt. Anytime bottom-up rule learning for knowledge graph completion. In *Proceedings of the 28th International Joint Conference on Artificial Intelligence, IJCAI'19*, page 3137–3143. AAAI Press, 2019.
 - 99 Saiping Guan, Jiyao Wei, Xiaolong Jin, Jiafeng Guo, and Xueqi Cheng. Look globally and reason: Two-stage path reasoning over sparse knowledge graphs. In *Proceedings of the 33rd ACM International Conference on Information and Knowledge Management, CIKM '24*, page 695–705, New York, NY, USA, 2024. Association for Computing Machinery.
 - 100 Wenhan Xiong, Thien Hoang, and William Yang Wang. DeepPath: A reinforcement learning method for knowledge graph reasoning. In Martha Palmer, Rebecca Hwa, and Sebastian Riedel, editors, *Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing*, pages 564–573, Copenhagen, Denmark, September 2017. Association for Computational Linguistics.
 - 101 Rajarshi Das, Shehzaad Dhuliawala, Manzil Zaheer, Luke Vilnis, Ishan Durugkar, Akshay Krishnamurthy, Alex Smola, and Andrew McCallum. Go for a walk and arrive at the answer: Reasoning over paths in knowledge bases using reinforcement learning. In *6th International Conference on Learning Representations, ICLR 2018, Vancouver, BC, Canada, April 30 – May 3, 2018, Conference Track Proceedings*. OpenReview.net, 2018.
 - 102 Xi Victoria Lin, Richard Socher, and Caiming Xiong. Multi-hop knowledge graph reasoning with reward shaping. In Ellen Riloff, David Chiang, Julia Hockenmaier, and Jun'ichi Tsujii, editors, *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 3243–3253, Brussels, Belgium, October–November 2018. Association for Computational Linguistics.
 - 103 Xin Lv, Xu Han, Lei Hou, Juanzi Li, Zhiyuan Liu, Wei Zhang, Yichi Zhang, Hao Kong, and Suhui Wu. Dynamic anticipation and completion for multi-hop reasoning over sparse knowledge graph. In Bonnie Webber, Trevor Cohn, Yulan He, and Yang Liu, editors, *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 5694–5703, Online, November 2020. Association for Computational Linguistics.
 - 104 Cheng Yan, Feng Zhao, and Hai Jin. Exkgr: Explainable multi-hop reasoning for evolving knowledge graph. In Arnab Bhattacharya, Janice Lee, Mong Li, Divyakant Agrawal, P. Krishna Reddy, Mukesh K. Mohania, Anirban Mondal, Vikram Goyal, and Rage Uday Kiran, editors, *Database Systems for Advanced Applications – 27th International Conference, DASFAA 2022, Virtual Event, April 11–14, 2022, Proceedings, Part I*, volume 13245 of *Lecture Notes in Computer Science*, pages 153–161. Springer, 2022.
 - 105 Zhongni Hou, Xiaolong Jin, Zixuan Li, and Long Bai. Rule-aware reinforcement learning for knowledge graph reasoning. In Chengqing Zong, Fei Xia, Wenjie Li, and Roberto Navigli, editors, *Findings of the Association for Computational Linguistics: ACL-IJCNLP 2021*, pages 4687–4692, Online, August 2021. Association for Computational Linguistics.
 - 106 Deren Lei, Gangrong Jiang, Xiaotao Gu, Kexuan Sun, Yuning Mao, and Xiang Ren. Learning collaborative agents with rule guidance for knowledge graph reasoning. In Bonnie Webber, Trevor Cohn, Yulan He, and Yang Liu, editors, *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 8541–8547, Online, November 2020. Association for Computational Linguistics.
 - 107 Xin Lv, Yuxian Gu, Xu Han, Lei Hou, Juanzi Li, and Zhiyuan Liu. Adapting meta knowledge graph information for multi-hop reasoning over few-shot relations. In Kentaro Inui, Jing Jiang, Vincent Ng, and Xiaojun Wan, editors, *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 3376–3381, Hong Kong, China, November 2019. Association for Computational Linguistics.
 - 108 Chuxu Zhang, Lu Yu, Mandana Saebi, Meng Jiang, and Nitesh Chawla. Few-shot multi-hop relation reasoning over knowledge bases. In Trevor Cohn, Yulan He, and Yang Liu, editors, *Findings of the Association for Computational Linguistics: EMNLP 2020*, pages 580–585, Online, November 2020. Association for Computational Linguistics.
 - 109 Shangfei Zheng, Wei Chen, Pengpeng Zhao, An Liu, Junhua Fang, and Lei Zhao. When hardness makes a difference: Multi-hop knowledge graph reasoning over few-shot relations. In *Proceedings of the 30th ACM International Conference on Information & Knowledge Management, CIKM '21*, page 2688–2697, New York, NY, USA, 2021. Association for Computing Machinery.
 - 110 Guojia Wan and Bo Du. Gaussianpath: A bayesian multi-hop reasoning framework for knowledge graph reasoning. In *AAAI Conference on Artificial Intelligence*, 2021.
 - 111 Guojia Wan, Shirui Pan, Chen Gong, Chuan Zhou, and Gholamreza Haffari. Reasoning like human: hierarchical reinforcement learning for knowledge graph reasoning. In *Proceedings of the Twenty-Ninth International Joint Conference on Artificial Intelligence, IJCAI'20*, 2021.
 - 112 Yushi Bai, Xin Lv, Juanzi Li, Lei Hou, Yincen Qu, Zelin Dai, and Feiyu Xiong. SQUIRE: A sequence-to-sequence framework for multi-hop knowledge graph reasoning. In Yoav Goldberg, Zornitsa Kozareva, and Yue Zhang, editors, *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing, EMNLP 2022, Abu Dhabi, United Arab Emirates, December 7–11, 2022*, pages 1649–1662. Association for Computational Linguistics, 2022.
 - 113 Mayi Xu, Ke Sun, Yongqi Li, and Tiejun Qian. Cold-start multi-hop reasoning by hierarchical guidance and self-verification. In *Machine Learning and Knowledge Discovery in Databases: Research Track: European Conference, ECML PKDD 2023, Turin, Italy, September 18–22, 2023, Proceedings, Part II*, page 577–592, Berlin, Heidelberg, 2023. Springer-Verlag.
 - 114 William L. Hamilton, Payal Bajaj, Marinka Zitnik, Dan Jurafsky, and Jure Leskovec. Embedding logical queries on knowledge graphs. In Samy Bengio, Hanna M. Wallach, Hugo Larochelle, Kristen Grauman, Nicolò Cesa-Bianchi, and Roman Garnett, editors, *Advances in Neural Information Processing Systems 31: Annual Conference on Neural Information Processing Systems 2018, NeurIPS 2018, December 3–8, 2018, Montréal, Canada*, pages 2030–2041, 2018.
 - 115 Hongyu Ren, Weihua Hu, and Jure Leskovec. Query2box: Reasoning over knowledge graphs in vector space using box embeddings. In *8th International Conference on Learning Representations, ICLR 2020, Addis Ababa, Ethiopia, April 26–30, 2020*. OpenReview.net, 2020.
 - 116 Nurendra Choudhary, Nikhil Rao, Sumeet Katariya, Karthik Subbian, and Chandan K. Reddy. Self-supervised hyperboloid representations from logical queries over knowledge graphs. In *Proceedings of the Web Conference 2021, WWW '21*, page 1373–1384, New York, NY, USA, 2021. Association for Computing Machinery.
 - 117 Zhanqiu Zhang, Jie Wang, Jiajun Chen, Shuiwang Ji, and Feng Wu. Cone: cone embeddings for multi-hop reasoning over knowledge graphs. In *Proceedings of the 35th International Conference on Neural Information Processing Systems, NIPS '21*, Red Hook, NY, USA, 2024. Curran Associates Inc.
 - 118 Hongyu Ren and Jure Leskovec. Beta embeddings for multi-hop logical reasoning in knowledge graphs. In *Proceedings of the 34th International Conference on Neural Information Processing Systems, NIPS '20*, Red Hook, NY, USA, 2020. Curran Associates Inc.
 - 119 Dong Yang, Peijun Qing, Yang Li, Haonan Lu, and Xiaodong Lin. GammaE: Gamma embeddings for logical queries on

- knowledge graphs. In Yoav Goldberg, Zornitsa Kozareva, and Yue Zhang, editors, *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pages 745–760, Abu Dhabi, United Arab Emirates, December 2022. Association for Computational Linguistics.
- 120 Alfonso Amayuelas, Shuai Zhang, Susie Xi Rao, and Ce Zhang. Neural methods for logical reasoning over knowledge graphs. In *The Tenth International Conference on Learning Representations (ICLR 2022)*. OpenReview, 2022.
 - 121 Haitian Sun, Andrew O. Arnold, Tania Bedrax-Weiss, Fernando Pereira, and William W. Cohen. Faithful embeddings for knowledge base queries. In *Proceedings of the 34th International Conference on Neural Information Processing Systems*, NIPS '20, Red Hook, NY, USA, 2020. Curran Associates Inc.
 - 122 Erik Arakelyan, Daniel Daza, Pasquale Minervini, and Michael Cochez. Complex query answering with neural link predictors. In *International Conference on Learning Representations*, 2021.
 - 123 Zhaocheng Zhu, Mikhail Galkin, Zuobai Zhang, and Jian Tang. Neural-symbolic models for logical queries on knowledge graphs. In Kamalika Chaudhuri, Stefanie Jegelka, Le Song, Csaba Szepesvari, Gang Niu, and Sivan Sabato, editors, *Proceedings of the 39th International Conference on Machine Learning*, volume 162 of *Proceedings of Machine Learning Research*, pages 27454–27478. PMLR, 17–23 Jul 2022.
 - 124 Xuelu Chen, Ziniu Hu, and Yizhou Sun. Fuzzy logic based logical query answering on knowledge graphs. *Proceedings of the AAAI Conference on Artificial Intelligence*, 36:3939–3948, 06 2022.
 - 125 Dingmin Wang, Yeyuan Chen, and Bernardo Cuenca Grau. Efficient embeddings of logical variables for query answering over incomplete knowledge graphs. *Proceedings of the AAAI Conference on Artificial Intelligence*, 37(4):4652–4659, Jun. 2023.
 - 126 Zhiwei Hu, Victor Gutierrez Basulto, Zhiliang Xiang, Xiaoli Li, Ru Li, and Jeff Z. Pan. Type-aware embeddings for multi-hop reasoning over knowledge graphs. In Lud De Raedt, editor, *Proceedings of the Thirty-First International Joint Conference on Artificial Intelligence, IJCAI-22*, pages 3078–3084. International Joint Conferences on Artificial Intelligence Organization, 7 2022. Main Track.
 - 127 Runyu Ni, Zhonggui Ma, Kaihang Yu, and Xiaohan Xu. Specific time embedding for temporal knowledge graph completion. In *19th IEEE International Conference on Cognitive Informatics & Cognitive Computing, ICCI*CC 2020, Beijing, China, September 26-28, 2020*, pages 105–110. IEEE, 2020.
 - 128 Lifan Lin and Kun She. Tensor decomposition-based temporal knowledge graph embedding. In *32nd IEEE International Conference on Tools with Artificial Intelligence, ICTAI 2020, Baltimore, MD, USA, November 9-11, 2020*, pages 969–975. IEEE, 2020.
 - 129 Johannes Messner, Ralph Abboud, and İsmail İlkan Ceylan. Temporal knowledge graph completion using box embeddings. In *Thirty-Sixth AAAI Conference on Artificial Intelligence, AAAI 2022, Thirty-Fourth Conference on Innovative Applications of Artificial Intelligence, IAAI 2022, The Twelveth Symposium on Educational Advances in Artificial Intelligence, EAAI 2022 Virtual Event, February 22 - March 1, 2022*, pages 7779–7787. AAAI Press, 2022.
 - 130 Ali Sadeghian, Mohammadreza Armandpour, Anthony M. Colas, and Daisy Zhe Wang. Chronor: Rotation based temporal knowledge graph embedding. In *Thirty-Fifth AAAI Conference on Artificial Intelligence, AAAI 2021, Thirty-Third Conference on Innovative Applications of Artificial Intelligence, IAAI 2021, The Eleventh Symposium on Educational Advances in Artificial Intelligence, EAAI 2021, Virtual Event, February 2-9, 2021*, pages 6471–6479. AAAI Press, 2021.
 - 131 Timothée Lacroix, Guillaume Obozinski, and Nicolas Usunier. Tensor decompositions for temporal knowledge base completion. In *8th International Conference on Learning Representations, ICLR 2020, Addis Ababa, Ethiopia, April 26-30, 2020*. OpenReview.net, 2020.
 - 132 Chengjin Xu, Mojtaba Nayyeri, Fouad Alkhoury, Hamed Shariat Yazdi, and Jens Lehmann. Tero: A time-aware knowledge graph embedding via temporal rotation. In Donia Scott, Núria Bel, and Chengqing Zong, editors, *Proceedings of the 28th International Conference on Computational Linguistics, COLING 2020, Barcelona, Spain (Online), December 8-13, 2020*, pages 1583–1593. International Committee on Computational Linguistics, 2020.
 - 133 Alberto García-Durán, Sebastijan Dumančić, and Mathias Niepert. Learning sequence encoders for temporal knowledge graph completion. In Ellen Riloff, David Chiang, Julia Hockenmaier, and Jun'ichi Tsujii, editors, *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 4816–4821, Brussels, Belgium, October-November 2018. Association for Computational Linguistics.
 - 134 Shib Sankar Dasgupta, Swayambhu Nath Ray, and Partha Talukdar. HyTE: Hyperplane-based temporally aware knowledge graph embedding. In Ellen Riloff, David Chiang, Julia Hockenmaier, and Jun'ichi Tsujii, editors, *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 2001–2011, Brussels, Belgium, October-November 2018. Association for Computational Linguistics.
 - 135 Rishab Goel, Seyed Mehran Kazemi, Marcus Brubaker, and Pascal Poupart. Diachronic embedding for temporal knowledge graph completion. In *Proceedings of the AAAI conference on artificial intelligence*, volume 34, pages 3988–3995, 2020.
 - 136 Prachi Jain, Sushant Rathi, Mausam, and Soumen Chakrabarti. Temporal Knowledge Base Completion: New Algorithms and Evaluation Protocols. In Bonnie Webber, Trevor Cohn, Yulan He, and Yang Liu, editors, *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 3733–3747, Online, November 2020. Association for Computational Linguistics.
 - 137 Luyi Bai, Wenting Yu, Mingzhuo Chen, and Xiangnan Ma. Multi-hop reasoning over paths in temporal knowledge graphs using reinforcement learning. *Appl. Soft Comput.*, 103:107144, 2021.
 - 138 Ye Tao, Ying Li, and Zhonghai Wu. Temporal link prediction via reinforcement learning. In *ICASSP 2021 - 2021 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 3470–3474, 2021.
 - 139 Haoai Sun, Jialun Zhong, Yunpu Ma, Zhen Han, and Kun He. TimeTraveler: Reinforcement learning for temporal knowledge graph forecasting. In Marie-Francine Moens, Xuanjing Huang, Lucia Specia, and Scott Wen-tau Yih, editors, *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 8306–8319, Online and Punta Cana, Dominican Republic, November 2021. Association for Computational Linguistics.
 - 140 Zixuan Li, Xiaolong Jin, Saiping Guan, Wei Li, Jiafeng Guo, Yuanzhuo Wang, and Xueqi Cheng. Search from history and reason for future: Two-stage reasoning on temporal knowledge graphs. In Chengqing Zong, Fei Xia, Wenjie Li, and Roberto Navigli, editors, *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pages 4732–4743, Online, August 2021. Association for Computational Linguistics.
 - 141 Yushan Liu, Yunpu Ma, Marcel Hildebrandt, Mitchell Joblin, and Volker Tresp. Tlogic: Temporal logical rules for explainable link forecasting on temporal knowledge graphs. In *Thirty-Sixth AAAI Conference on Artificial Intelligence, AAAI 2022, Thirty-Fourth Conference on Innovative Applications of Artificial Intelligence, IAAI 2022, The Twelveth Symposium on Educational Advances in Artificial Intelligence, EAAI 2022 Virtual Event, February 22 - March 1, 2022*, pages 4120–4127. AAAI Press, 2022.
 - 142 Guanglin Niu and Bo Li. Logic and commonsense-guided temporal knowledge graph completion. In Brian Williams, Yiling Chen, and Jennifer Neville, editors, *Thirty-Seventh AAAI Conference on Artificial Intelligence, AAAI 2023, Thirty-Fifth Conference on Innovative Applications of Artificial Intelligence, IAAI 2023, Thirteenth Symposium on Educational Advances in Artificial Intelligence, EAAI 2023, Washington, DC, USA, February 7-14, 2023*, pages 4569–4577. AAAI Press, 2023.

- 143 Luyi Bai, Mingzhuo Chen, Lin Zhu, and Xiangxi Meng. Multi-hop temporal knowledge graph reasoning with temporal path rules guidance. *Expert Syst. Appl.*, 223:119804, 2023.
- 144 Shangfei Zheng, Hongzhi Yin, Tong Chen, Quoc Viet Hung Nguyen, Wei Chen, and Lei Zhao. Dream: Adaptive reinforcement learning based on attention mechanism for temporal knowledge graph reasoning. In *Proceedings of the 46th International ACM SIGIR Conference on Research and Development in Information Retrieval*, SIGIR '23, page 1578–1588, New York, NY, USA, 2023. Association for Computing Machinery.
- 145 Rakshit Trivedi, Hanjun Dai, Yichen Wang, and Le Song. Know-evolve: deep temporal reasoning for dynamic knowledge graphs. In *Proceedings of the 34th International Conference on Machine Learning - Volume 70*, ICML'17, page 3462–3471. JMLR.org, 2017.
- 146 Woojeong Jin, Meng Qu, Xisen Jin, and Xiang Ren. Recurrent event network: Autoregressive structure inference over temporal knowledge graphs. In Bonnie Webber, Trevor Cohn, Yulan He, and Yang Liu, editors, *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 6669–6683, Online, November 2020. Association for Computational Linguistics.
- 147 Zixuan Li, Xiaolong Jin, Wei Li, Saiping Guan, Jiafeng Guo, Huawei Shen, Yuanzhuo Wang, and Xueqi Cheng. Temporal knowledge graph reasoning based on evolutionary representation learning. In Fernando Diaz, Chirag Shah, Torsten Suel, Pablo Castells, Rosie Jones, and Tetsuya Sakai, editors, *SIGIR '21: The 44th International ACM SIGIR Conference on Research and Development in Information Retrieval, Virtual Event, Canada, July 11–15, 2021*, pages 408–417. ACM, 2021.
- 148 Zixuan Li, Saiping Guan, Xiaolong Jin, Weihua Peng, Yajuan Lyu, Yong Zhu, Long Bai, Wei Li, Jiafeng Guo, and Xueqi Cheng. Complex evolutionary pattern learning for temporal knowledge graph reasoning. In Smaranda Muresan, Preslav Nakov, and Aline Villavicencio, editors, *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers), ACL 2022, Dublin, Ireland, May 22–27, 2022*, pages 290–296. Association for Computational Linguistics, 2022.
- 149 Cunchao Zhu, Muhao Chen, Changjun Fan, Guangquan Cheng, and Yan Zhang. Learning from history: Modeling temporal knowledge graphs with sequential copy-generation networks. In *Thirty-Fifth AAAI Conference on Artificial Intelligence, AAAI 2021, Thirty-Third Conference on Innovative Applications of Artificial Intelligence, IAAI 2021, The Eleventh Symposium on Educational Advances in Artificial Intelligence, EAAI 2021, Virtual Event, February 2–9, 2021*, pages 4732–4740. AAAI Press, 2021.
- 150 Yi Xu, Junjie Ou, Hui Xu, and Luoyi Fu. Temporal knowledge graph reasoning with historical contrastive learning. In Brian Williams, Yiling Chen, and Jennifer Neville, editors, *Thirty-Seventh AAAI Conference on Artificial Intelligence, AAAI 2023, Thirty-Fifth Conference on Innovative Applications of Artificial Intelligence, IAAI 2023, Thirteenth Symposium on Educational Advances in Artificial Intelligence, EAAI 2023, Washington, DC, USA, February 7–14, 2023*, pages 4765–4773. AAAI Press, 2023.
- 151 Yujia Li, Shiliang Sun, and Jing Zhao. Tirgn: Time-guided recurrent graph network with local-global historical patterns for temporal knowledge graph reasoning. In Luc De Raedt, editor, *Proceedings of the Thirty-First International Joint Conference on Artificial Intelligence, IJCAI 2022, Vienna, Austria, 23–29 July 2022*, pages 2152–2158. ijcai.org, 2022.
- 152 Mengqi Zhang, Yuwei Xia, Qiang Liu, Shu Wu, and Liang Wang. Learning long- and short-term representations for temporal knowledge graph reasoning. In Ying Ding, Jie Tang, Juan F. Sequeda, Lora Aroyo, Carlos Castillo, and Geert-Jan Houben, editors, *Proceedings of the ACM Web Conference 2023, WWW 2023, Austin, TX, USA, 30 April 2023 - 4 May 2023*, pages 2412–2422. ACM, 2023.
- 153 Ruobing Xie, Zhiyuan Liu, Huanbo Luan, and Maosong Sun. Image-embodied knowledge representation learning. In *Proceedings of the 26th International Joint Conference on Artificial Intelligence, IJCAI'17*, page 3140–3146. AAAI Press, 2017.
- 154 Hatem Mousselly-Sergieh, Teresa Botschen, Iryna Gurevych, and Stefan Roth. A multimodal translation-based approach for knowledge graph representation learning. In Malvina Nissim, Jonathan Berant, and Alessandro Lenci, editors, *Proceedings of the Seventh Joint Conference on Lexical and Computational Semantics*, pages 225–234, New Orleans, Louisiana, June 2018. Association for Computational Linguistics.
- 155 Zikang Wang, Linjing Li, Qiudan Li, and Daniel Zeng. Multimodal data enhanced representation learning for knowledge graphs. In *2019 International Joint Conference on Neural Networks (IJCNN)*, pages 1–8, 2019.
- 156 Pouya Pezeshkpour, Liyan Chen, and Sameer Singh. Embedding multimodal relational data for knowledge base completion. In Ellen Riloff, David Chiang, Julia Hockenmaier, and Jun'ichi Tsujii, editors, *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 3208–3218, Brussels, Belgium, October–November 2018. Association for Computational Linguistics.
- 157 Steffen Thoma, Achim Rettinger, and Fabian Both. Towards holistic concept representations: Embedding relational knowledge, visual attributes, and distributional word semantics. In *The Semantic Web – ISWC 2017: 16th International Semantic Web Conference, Vienna, Austria, October 21–25, 2017, Proceedings, Part I*, page 694–710, Berlin, Heidelberg, 2017. Springer-Verlag.
- 158 Xinyu Lu, Lifang Wang, Zejun Jiang, Shichang He, and Shizhong Liu. Mmkr1: A robust embedding approach for multi-modal knowledge graph representation learning. *Applied Intelligence*, 52(7):7480–7497, May 2022.
- 159 Meng Wang, Sen Wang, Han Yang, Zheng Zhang, Xi Chen, and Guilin Qi. Is visual context really helpful for knowledge graph? a representation learning perspective. In *Proceedings of the 29th ACM International Conference on Multimedia, MM '21*, page 2735–2743, New York, NY, USA, 2021. Association for Computing Machinery.
- 160 Yichi Zhang and Wen Zhang. Knowledge graph completion with pre-trained multimodal transformer and twins negative sampling, 2022.
- 161 Xiang Chen, Ningyu Zhang, Lei Li, Shumin Deng, Chuanqi Tan, Changliang Xu, Fei Huang, Luo Si, and Huajun Chen. Hybrid transformer with multi-level fusion for multimodal knowledge graph completion. In *Proceedings of the 45th International ACM SIGIR Conference on Research and Development in Information Retrieval*, SIGIR '22, page 904–915, New York, NY, USA, 2022. Association for Computing Machinery.
- 162 Ke Liang, Lingyuan Meng, Yue Liu, Meng Wei, Wei Wei, Suyuan Liu, Wenxuan Tu, Siwei Wang, Sihang Zhou, and Xinwang Liu. Simple yet effective: Structure guided pre-trained transformer for multi-modal knowledge graph reasoning. In *Proceedings of the 32nd ACM International Conference on Multimedia, MM '24*, page 1554–1563, New York, NY, USA, 2024. Association for Computing Machinery.
- 163 Mehdi Ali, Max Berrendorf, Charles Tapley Hoyt, Laurent Vermue, Mikhail Galkin, Sahand Sharifzadeh, Asja Fischer, Volker Tresp, and Jens Lehmann. Bringing light into the dark: A large-scale evaluation of knowledge graph embedding models under a unified framework. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 44(12):8825–8845, 2022.
- 164 Jie Zhou, Ganqu Cui, Shengding Hu, Zhengyan Zhang, Cheng Yang, Zhiyuan Liu, Lifeng Wang, Changcheng Li, and Maosong Sun. Graph neural networks: A review of methods and applications. *AI Open*, 1:57–81, January 2020.
- 165 Chelsea Finn, Pieter Abbeel, and Sergey Levine. Model-agnostic meta-learning for fast adaptation of deep networks. In *Proceedings of the 34th International Conference on Machine Learning - Volume 70*, ICML'17, page 1126–1135. JMLR.org, 2017.
- 166 Ralph Abboud, İsmail İlkan Ceylan, Thomas Lukasiewicz, and Tommaso Salvatori. Boxe: A box embedding model for

- knowledge base completion. In Hugo Larochelle, Marc'Aurelio Ranzato, Raia Hadsell, Maria-Florina Balcan, and Hsuantien Lin, editors, *Advances in Neural Information Processing Systems 33: Annual Conference on Neural Information Processing Systems 2020, NeurIPS 2020, December 6-12, 2020, virtual*, 2020.
- 167 Thomas N. Kipf and Max Welling. Semi-supervised classification with graph convolutional networks. In *International Conference on Learning Representations*, 2017.
- 168 Liunian Harold Li, Mark Yatskar, Da Yin, Cho-Jui Hsieh, and Kai-Wei Chang. Visualbert: A simple and performant baseline for vision and language. *ArXiv*, abs/1908.03557, 2019.
- 169 Maarten Sap, Ronan Le Bras, Emily Allaway, Chandra Bhagavatula, Nicholas Lourie, Hannah Rashkin, Brendan Roof, Noah A. Smith, and Yejin Choi. Atomic: An atlas of machine commonsense for if-then reasoning. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 33, pages 3027–3035, 2019.
- 170 Tara Safavi and Danai Koutra. Codex: A compositional and domain-aware benchmark for knowledge graph completion. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, 2020.
- 171 Robert Speer, Joshua Chin, and Catherine Havasi. Conceptnet 5.5: An open multilingual graph of general knowledge. *Proceedings of the AAAI Conference on Artificial Intelligence*, 31(1):4444–4451, 2017.
- 172 Xiang Li, Aynaz Taheri, Lifu Tu, and Kevin Gimpel. Commonsense knowledge base completion. In Katrin Erk and Noah A. Smith, editors, *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1445–1455, Berlin, Germany, August 2016. Association for Computational Linguistics.
- 173 Baoxu Shi and Tim Wengener. Open-world knowledge graph completion. In *Proceedings of the AAAI conference on artificial intelligence*, volume 32, 2018.
- 174 Kristina Toutanova, Danqi Chen, Patrick Pantel, Hoifung Poon, Pallavi Choudhury, and Michael Gamon. Observed versus latent features for knowledge base and text inference. In *Proceedings of the 3rd Workshop on Continuous Vector Space Models and their Compositionality*, pages 57–66, 2015.
- 175 Daniel Scott Himmelstein, Antoine Lizee, Christine Hessler, Leo Brueggeman, Sabrina L Chen, Dexter Hadley, Ari Green, Pouya Khankhanian, and Sergio E Baranzini. Systematic integration of biomedical knowledge prioritizes drugs for repurposing. *eLife*, 6:e26726, sep 2017.
- 176 Anna Breit, Simon Ott, Asan Agibetov, and Matthias Samwald. Openbiolink: a benchmarking framework for large-scale biomedical link prediction. *Bioinformatics*, 36(13):4097–4098, 04 2020.
- 177 Charles Kemp, Joshua B. Tenenbaum, Thomas L. Griffiths, Takeshi Yamada, and Naonori Ueda. Learning systems of concepts with an infinite relational model. In *Proceedings of the 21st National Conference on Artificial Intelligence - Volume 1, AAAI'06*, page 381–388. AAAI Press, 2006.
- 178 Tim Dettmers, Pasquale Minervini, Pontus Stenetorp, and Sebastian Riedel. Convolutional 2d knowledge graph embeddings. In *Proceedings of the Thirty-Second AAAI Conference on Artificial Intelligence and Thirtieth Innovative Applications of Artificial Intelligence Conference and Eighth AAAI Symposium on Educational Advances in Artificial Intelligence, AAAI'18/IAAI'18/EAAI'18*. AAAI Press, 2018.
- 179 Tianxing Wu, Arijit Khan, Melvin Yong, Guilin Qi, and Meng Wang. Efficiently embedding dynamic knowledge graphs. *Knowledge-Based Systems*, 250:109124, 2022.
- 180 Fuwei Zhang, Zhao Zhang, Xiang Ao, Fuzhen Zhuang, Yongjun Xu, and Qing He. Along the time: Timeline-traced embedding for temporal knowledge graph completion. In *Proceedings of the 31st ACM International Conference on Information & Knowledge Management, CIKM '22*, page 2529–2538, New York, NY, USA, 2022. Association for Computing Machinery.
- 181 Sebastián Ferrada, Benjamin Bustos, and Aidan Hogan. Imgpedia: A linked dataset with content-based analysis of wikipedia images. In Claudia d'Amato, Miriam Fernandez, Valentina Tamma, Freddy Lecue, Philippe Cudré-Mauroux, Juan Sequeda, Christoph Lange, and Jeff Hefflin, editors, *The Semantic Web – ISWC 2017*, pages 84–93, Cham, 2017. Springer International Publishing.
- 182 Ye Liu, Hui Li, Alberto Garcia-Duran, Mathias Niepert, Daniel Onoro-Rubio, and David S. Rosenblum. Mmkg: Multi-modal knowledge graphs. In Pascal Hitzler, Miriam Fernández, Krzysztof Janowicz, Amrapali Zaveri, Alasdair J.G. Gray, Vanessa Lopez, Armin Haller, and Karl Hammar, editors, *The Semantic Web*, pages 459–474, Cham, 2019. Springer International Publishing.
- 183 Derong Xu, Tong Xu, Shiwei Wu, Jingbo Zhou, and Enhong Chen. Relation-enhanced negative sampling for multimodal knowledge graph completion. In *Proceedings of the 30th ACM International Conference on Multimedia, MM '22*, page 3857–3866, New York, NY, USA, 2022. Association for Computing Machinery.
- 184 Meng Wang, Haofen Wang, Guilin Qi, and Qiushuo Zheng. Richpedia: A large-scale, comprehensive multi-modal knowledge graph. *Big Data Research*, 22:100159, 2020.
- 185 Nazneen Fatema Rajani, Bryan McCann, Caiming Xiong, and Richard Socher. Explain yourself! leveraging language models for commonsense reasoning. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 4932–4942, 2019.
- 186 Alon Talmor, Oyvind Tafjord, Peter Clark, Yoav Goldberg, and Jonathan Berant. Leap-of-thought: Teaching pre-trained models to systematically reason over implicit knowledge. *Advances in Neural Information Processing Systems*, 33:20227–20237, 2020.
- 187 Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, Ilya Sutskever, et al. Language models are unsupervised multitask learners. 2019.
- 188 Charles R Fletcher. Understanding and solving arithmetic word problems: A computer simulation. *Behavior Research Methods, Instruments, & Computers*, 17(5):565–571, 1985.
- 189 Ma Yuhui, Zhou Ying, Cui Guangzuo, Ren Yun, and Huang Ronghuai. Frame-based calculus of solving arithmetic multi-step addition and subtraction word problems. In *2010 Second International Workshop on Education Technology and Computer Science*, volume 2, page 476–479. IEEE, 2010.
- 190 Christian Liguda and Thies Pfeiffer. Modeling math word problems with augmented semantic networks. In *International Conference on Application of Natural Language to Information Systems*, page 247–252. Springer, 2012.
- 191 Shuming Shi, Yuehui Wang, Chin-Yew Lin, Xiaojiang Liu, and Yong Rui. Automatically solving number word problems by semantic parsing and reasoning. In *Proceedings of the 2015 conference on empirical methods in natural language processing*, page 1132–1142, 2015.
- 192 Rik Koncel-Kedziorski, Hannaneh Hajishirzi, Ashish Sabharwal, Oren Etzioni, and Siena Dumas Ang. Parsing algebraic word problems into equations. *Transactions of the Association for Computational Linguistics*, 3:585–597, 2015.
- 193 Shyam Upadhyay, Ming-Wei Chang, Kai-Wei Chang, and Wen-tau Yih. Learning from explicit and implicit supervision jointly for algebra word problems. In *Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing*, page 297–306, 2016.
- 194 Yao Fu, Hao Peng, Ashish Sabharwal, Peter Clark, and Tushar Khot. Complexity-based prompting for multi-step reasoning. In *The Eleventh International Conference on Learning Representations, ICLR 2023, Kigali, Rwanda, May 1-5, 2023*. OpenReview.net, 2023.
- 195 Yew Ken Chia, Guizhen Chen, Luu Anh Tuan, Soujanya Poria, and Lidong Bing. Contrastive chain-of-thought prompting. *arXiv preprint arXiv:2311.09277*, 2023.
- 196 Shizhe Diao, Pengcheng Wang, LIN Yong, Rui Pan, Xiang Liu, and Tong Zhang. Active prompting with chain-of-thought

- for large language models, 2024.
- 197 Mayi Xu, Yongqi Li, Ke Sun, and Tiejun Qian. Adaption-of-thought: Learning question difficulty improves large language models for reasoning. In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, pages 5468–5495, 2024.
- 198 Shunyu Yao, Dian Yu, Jeffrey Zhao, Izhak Shafran, Tom Griffiths, Yuan Cao, and Karthik Narasimhan. Tree of thoughts: Deliberate problem solving with large language models. *Advances in Neural Information Processing Systems*, 36, 2024.
- 199 Maciej Besta, Nils Blach, Ales Kubicek, Robert Gerstenberger, Michal Podstawski, Lukas Gianinazzi, Joanna Gajda, Tomasz Lehmann, Hubert Niewiadomski, Piotr Nyczyk, et al. Graph of thoughts: Solving elaborate problems with large language models. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 38, pages 17682–17690, 2024.
- 200 Zhiyuan Hu, Chumin Liu, Xidong Feng, Yilun Zhao, See-Kiong Ng, Anh Tuan Luu, Junxian He, Pang Wei Koh, and Bryan Hooi. Uncertainty of thoughts: Uncertainty-aware planning enhances information seeking in large language models. In *ICLR 2024 Workshop on Large Language Model (LLM) Agents*, 2024.
- 201 Chuanyang Zheng, Zhengying Liu, Enze Xie, Zhenguo Li, and Yu Li. Progressive-hint prompting improves reasoning in large language models. *arXiv preprint arXiv:2304.09797*, 2023.
- 202 Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are few-shot learners. *Advances in neural information processing systems*, 33:1877–1901, 2020.
- 203 Takeshi Kojima, Shixiang Shane Gu, Machel Reid, Yutaka Matsuo, and Yusuke Iwasawa. Large language models are zero-shot reasoners. *Advances in neural information processing systems*, 35:22199–22213, 2022.
- 204 Xinxin Lyu, Sewon Min, Iz Beltagy, Luke Zettlemoyer, and Hannaneh Hajishirzi. Z-icl: Zero-shot in-context learning with pseudo-demonstrations. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 2304–2317, 2023.
- 205 Zhuosheng Zhang, Aston Zhang, Mu Li, and Alex Smola. Automatic chain of thought prompting in large language models. In *The Eleventh International Conference on Learning Representations, ICLR 2023, Kigali, Rwanda, May 1-5, 2023*. OpenReview.net, 2023.
- 206 Xingchen Wan, Ruoxi Sun, Hanjun Dai, Sercan Arik, and Tomas Pfister. Better zero-shot reasoning with self-adaptive prompting. In *Findings of the Association for Computational Linguistics: ACL 2023*, pages 3493–3514, 2023.
- 207 Xiaonan Li and Xipeng Qiu. Mot: Memory-of-thought enables chatgpt to self-improve. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 6354–6374, 2023.
- 208 Noah Shinn, Federico Cassano, Ashwin Gopinath, Karthik Narasimhan, and Shunyu Yao. Reflexion: language agents with verbal reinforcement learning. In Alice Oh, Tristan Naumann, Amir Globerson, Kate Saenko, Moritz Hardt, and Sergey Levine, editors, *Advances in Neural Information Processing Systems 36: Annual Conference on Neural Information Processing Systems 2023, NeurIPS 2023, New Orleans, LA, USA, December 10 - 16, 2023*, 2023.
- 209 Yixuan Weng, Minjun Zhu, Fei Xia, Bin Li, Shizhu He, Shengping Liu, Bin Sun, Kang Liu, and Jun Zhao. Large language models are better reasoners with self-verification. In Houda Bouamor, Juan Pino, and Kalika Bali, editors, *Findings of the Association for Computational Linguistics: EMNLP 2023, Singapore, December 6-10, 2023*, pages 2550–2575. Association for Computational Linguistics, 2023.
- 210 Deepak Nathani, David Wang, Liangming Pan, and William Yang Wang. MAF: multi-aspect feedback for improving reasoning in large language models. In Houda Bouamor, Juan Pino, and Kalika Bali, editors, *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing, EMNLP 2023, Singapore, December 6-10, 2023*, pages 6591–6616. Association for Computational Linguistics, 2023.
- 211 Zhibin Gou, Zhihong Shao, Yeyun Gong, Yelong Shen, Yujiu Yang, Nan Duan, and Weizhu Chen. CRITIC: large language models can self-correct with tool-interactive critiquing. In *The Twelfth International Conference on Learning Representations, ICLR 2024, Vienna, Austria, May 7-11, 2024*. OpenReview.net, 2024.
- 212 Liangming Pan, Alon Albalak, Xinyi Wang, and William Yang Wang. Logic-lm: Empowering large language models with symbolic solvers for faithful logical reasoning. In Houda Bouamor, Juan Pino, and Kalika Bali, editors, *Findings of the Association for Computational Linguistics: EMNLP 2023, Singapore, December 6-10, 2023*, pages 3806–3824. Association for Computational Linguistics, 2023.
- 213 Wenhu Chen, Xueguang Ma, Xinyi Wang, and William W. Cohen. Program of thoughts prompting: Disentangling computation from reasoning for numerical reasoning tasks. *Trans. Mach. Learn. Res.*, 2023, 2023.
- 214 Luyu Gao, Aman Madaan, Shuyan Zhou, Uri Alon, Pengfei Liu, Yiming Yang, Jamie Callan, and Graham Neubig. PAL: Program-aided language models. In Andreas Krause, Emma Brunskill, Kyunghyun Cho, Barbara Engelhardt, Sivan Sabato, and Jonathan Scarlett, editors, *Proceedings of the 40th International Conference on Machine Learning, volume 202 of Proceedings of Machine Learning Research*, pages 10764–10799. PMLR, 23–29 Jul 2023.
- 215 Pan Lu, Baolin Peng, Hao Cheng, Michel Galley, Kai-Wei Chang, Ying Nian Wu, Song-Chun Zhu, and Jianfeng Gao. Chameleon: Plug-and-play compositional reasoning with large language models. In Alice Oh, Tristan Naumann, Amir Globerson, Kate Saenko, Moritz Hardt, and Sergey Levine, editors, *Advances in Neural Information Processing Systems 36: Annual Conference on Neural Information Processing Systems 2023, NeurIPS 2023, New Orleans, LA, USA, December 10 - 16, 2023*, 2023.
- 216 Ansong Ni, Sridi Iyer, Dragomir Radev, Veselin Stoyanov, Wen-Tau Yih, Sida Wang, and Xi Victoria Lin. LEVER: Learning to verify language-to-code generation with execution. In Andreas Krause, Emma Brunskill, Kyunghyun Cho, Barbara Engelhardt, Sivan Sabato, and Jonathan Scarlett, editors, *Proceedings of the 40th International Conference on Machine Learning*, volume 202 of *Proceedings of Machine Learning Research*, pages 26106–26128. PMLR, 23–29 Jul 2023.
- 217 Jonathan Uesato, Nate Kushman, Ramana Kumar, H. Francis Song, Noah Y. Siegel, Lisa Wang, Antonia Creswell, Geoffrey Irving, and Irina Higgins. Solving math word problems with process- and outcome-based feedback. *CoRR*, abs/2211.14275, 2022.
- 218 Hunter Lightman, Vineet Kosaraju, Yuri Burda, Harrison Edwards, Bowen Baker, Teddy Lee, Jan Leike, John Schulman, Ilya Sutskever, and Karl Cobbe. Let’s verify step by step. In *The Twelfth International Conference on Learning Representations, ICLR 2024, Vienna, Austria, May 7-11, 2024*. OpenReview.net, 2024.
- 219 Muhammad Khalifa, Lajanugen Logeswaran, Moontae Lee, Honglak Lee, and Lu Wang. GRACE: discriminator-guided chain-of-thought reasoning. In Houda Bouamor, Juan Pino, and Kalika Bali, editors, *Findings of the Association for Computational Linguistics: EMNLP 2023, Singapore, December 6-10, 2023*, pages 15299–15328. Association for Computational Linguistics, 2023.
- 220 Yifei Li, Zeqi Lin, Shizhuo Zhang, Qiang Fu, Bei Chen, Jian-Guang Lou, and Weizhu Chen. Making language models better reasoners with step-aware verifier. In Anna Rogers, Jordan L. Boyd-Graber, and Naoaki Okazaki, editors, *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), ACL 2023, Toronto, Canada, July 9-14, 2023*, pages 5315–5333. Association for Computational Linguistics, 2023.
- 221 Danqing Wang and Lei Li. Learning from mistakes via cooperative study assistant for large language models. In Houda Bouamor, Juan Pino, and Kalika Bali, editors, *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing, EMNLP 2023, Singapore, December 6-10, 2023*, pages 10667–10685. Association for Computational Linguistics, 2023.

- 222 Debjit Paul, Mete Ismayilzada, Maxime Peyrard, Beatriz Borges, Antoine Bosselut, Robert West, and Boi Faltings. RE-FINER: reasoning feedback on intermediate representations. In Yvette Graham and Matthew Purver, editors, *Proceedings of the 18th Conference of the European Chapter of the Association for Computational Linguistics, EACL 2024 - Volume 1: Long Papers, St. Julian's, Malta, March 17-22, 2024*, pages 1100–1126. Association for Computational Linguistics, 2024.
- 223 Ofir Press, Muru Zhang, Sewon Min, Ludwig Schmidt, Noah A. Smith, and Mike Lewis. Measuring and narrowing the compositionality gap in language models. In *Findings of the 2023 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, 2023.
- 224 Dheeru Dua, Shivanshu Gupta, Sameer Singh, and Matt Gardner. Successive prompting for decomposing complex questions. In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pages 1251–1265, 2022.
- 225 Lei Wang, Wanyu Xu, Yihuai Lan, Zhiqiang Hu, Yunshi Lan, Roy Ka-Wei Lee, and Ee-Peng Lim. Plan-and-solve prompting: Improving zero-shot chain-of-thought reasoning by large language models. In Anna Rogers, Jordan L. Boyd-Graber, and Naoaki Okazaki, editors, *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), ACL 2023, Toronto, Canada, July 9-14, 2023*, pages 2609–2634. Association for Computational Linguistics, 2023.
- 226 Denny Zhou, Nathanael Schärli, Le Hou, Jason Wei, Nathan Scales, Xuezhi Wang, Dale Schuurmans, Claire Cui, Olivier Bousquet, Quoc V. Le, and Ed H. Chi. Least-to-most prompting enables complex reasoning in large language models. In *The Eleventh International Conference on Learning Representations, ICLR 2023, Kigali, Rwanda, May 1-5, 2023*. OpenReview.net, 2023.
- 227 Andrew Drozdov, Nathanael Schärli, Ekin Akyürek, Nathan Scales, Xinying Song, Xinyun Chen, Olivier Bousquet, and Denny Zhou. Compositional semantic parsing with large language models. In *The Eleventh International Conference on Learning Representations*, 2022.
- 228 Tushar Khot, Harsh Trivedi, Matthew Finlayson, Yao Fu, Kyle Richardson, Peter Clark, and Ashish Sabharwal. Decomposed prompting: A modular approach for solving complex tasks. In *The Eleventh International Conference on Learning Representations*, 2023.
- 229 Xuefei Ning, Zinan Lin, Zixuan Zhou, Zifu Wang, Huazhong Yang, and Yu Wang. Skeleton-of-thought: Large language models can do parallel decoding. *Proceedings ENLSP-III*, 2023.
- 230 Yifei Li, Zeqi Lin, Shizhuo Zhang, Qiang Fu, Bei Chen, Jian-Guang Lou, and Weizhu Chen. Making large language models better reasoners with step-aware verifier. *arXiv preprint arXiv:2206.02336*, 2022.
- 231 Ori Yoran, Tomer Wolfson, Ben Bogin, Uri Katz, Daniel Deutch, and Jonathan Berant. Answering questions by meta-reasoning over multiple chains of thought. In *The 2023 Conference on Empirical Methods in Natural Language Processing*.
- 232 Justin Chen, Swarnadeep Saha, and Mohit Bansal. ReConcile: Round-table conference improves reasoning via consensus among diverse LLMs. In Lun-Wei Ku, Andre Martins, and Vivek Srikrumar, editors, *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 7066–7085, Bangkok, Thailand, August 2024. Association for Computational Linguistics.
- 233 Chi-Min Chan, Weize Chen, Yusheng Su, Jianxuan Yu, Wei Xue, Shanghang Zhang, Jie Fu, and Zhiyuan Liu. Chateval: Towards better llm-based evaluators through multi-agent debate. In *The Twelfth International Conference on Learning Representations*.
- 234 Luke Yoffe, Alfonso Amayuelas, and William Yang Wang. Debunc: Mitigating hallucinations in large language model agent communication with uncertainty estimations. *arXiv preprint arXiv:2407.06426*, 2024.
- 235 Qineng Wang, Zihao Wang, Ying Su, Hanghang Tong, and Yangqiu Song. Rethinking the bounds of LLM reasoning: Are multi-agent discussions the key? In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 6106–6131, 2024.
- 236 Weize Chen, Yusheng Su, Jingwei Zuo, Cheng Yang, Chenfei Yuan, Chi-Min Chan, Heyang Yu, Yaxi Lu, Yi-Hsin Hung, Chen Qian, et al. Agentverse: Facilitating multi-agent collaboration and exploring emergent behaviors. In *The Twelfth International Conference on Learning Representations*, 2023.
- 237 Kai Xiong, Xiao Ding, Yixin Cao, Ting Liu, and Bing Qin. Examining inter-consistency of large language models collaboration: An in-depth analysis via debate. In *Findings of the Association for Computational Linguistics: EMNLP 2023*, pages 7572–7590, 2023.
- 238 Bin Lei, Yi Zhang, Shan Zuo, Ali Payani, and Caiwen Ding. Macm: Utilizing a multi-agent system for condition mining in solving complex mathematical problems. *arXiv preprint arXiv:2404.04735*, 2024.
- 239 Xiaohe Bo, Zeyu Zhang, Quanyu Dai, Xueyang Feng, Lei Wang, Rui Li, Xu Chen, and Ji-Rong Wen. Reflective multi-agent collaboration based on large language models. In *The Thirty-eighth Annual Conference on Neural Information Processing Systems*.
- 240 Dawei Li, Zhen Tan, Peijia Qian, Yifan Li, Kumar Satvik Chaudhary, Lijie Hu, and Jiayi Shen. Smoa: Improving multi-agent large language models with sparse mixture-of-agents. *arXiv preprint arXiv:2411.03284*, 2024.
- 241 Liangming Pan, Michael Saxon, Wenda Xu, Deepak Nathani, Xinyi Wang, and William Yang Wang. Automatically correcting large language models: *Surveying the Landscape of Diverse Automated Correction Strategies*. *Trans. Assoc. Comput. Linguistics*, 12:484–506, 2024.
- 242 Jie Huang, Xinyun Chen, Swaroop Mishra, Huaixiu Steven Zheng, Adams Wei Yu, Xinying Song, and Denny Zhou. Large language models cannot self-correct reasoning yet. In *The Twelfth International Conference on Learning Representations, ICLR 2024, Vienna, Austria, May 7-11, 2024*. OpenReview.net, 2024.
- 243 Ryo Kamoi, Yusen Zhang, Nan Zhang, Jiawei Han, and Rui Zhang. When can llms actually correct their own mistakes? a critical survey of self-correction of llms. *Transactions of the Association for Computational Linguistics*, 12:1417–1440, 11 2024.
- 244 Marvin Minsky. *Society of mind*. Simon and Schuster, 1988.
- 245 Jintian Zhang, Xin Xu, Ningyu Zhang, Ruibo Liu, Bryan Hooi, and Shumin Deng. Exploring collaboration mechanisms for LLM agents: A social psychology view. In *ICLR 2024 Workshop on Large Language Model (LLM) Agents*, 2024.
- 246 Wang Ling, Dani Yogatama, Chris Dyer, and Phil Blunsom. Program induction by rationale generation: Learning to solve and explain algebraic word problems. In *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 158–167, 2017.
- 247 Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser, Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, et al. Training verifiers to solve math word problems. *arXiv preprint arXiv:2110.14168*, 2021.
- 248 Arkil Patel, Satwik Bhattamishra, and Navin Goyal. Are nlp models really able to solve simple math word problems? In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 2080–2094, 2021.
- 249 Mohammad Javad Hosseini, Hannaneh Hajishirzi, Oren Etzioni, and Nate Kushman. Learning to solve arithmetic word problems with verb categorization. In *Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, Jan 2014.
- 250 Subhro Roy and Dan Roth. Solving general arithmetic word problems. In *Proceedings of the 2015 Conference on Empirical*

- Methods in Natural Language Processing*. Association for Computational Linguistics, 2015.
- 251 Rik Koncel-Kedziorski, Hannaneh Hajishirzi, Ashish Sabharwal, Oren Etzioni, and Siena Dumas Ang. Parsing algebraic word problems into equations. *Transactions of the Association for Computational Linguistics*, page 585–597, Dec 2015.
- 252 Alon Talmor, Jonathan Herzig, Nicholas Lourie, and Jonathan Berant. Commonsenseqa: A question answering challenge targeting commonsense knowledge. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 4149–4158, 2019.
- 253 Mor Geva, Daniel Khashabi, Elad Segal, Tushar Khot, Dan Roth, and Jonathan Berant. Did aristotle use a laptop? a question answering benchmark with implicit reasoning strategies. *Transactions of the Association for Computational Linguistics*, 9:346–361, 2021.
- 254 Rajarshi Das, Manzil Zaheer, Dung Thai, Ameya Godbole, Ethan Perez, Jay Yoon Lee, Lizhen Tan, Lazaros Polymenakos, and Andrew McCallum. Case-based reasoning for natural language queries over knowledge bases. In Marie-Francine Moens, Xuanjing Huang, Lucia Specia, and Scott Wen-tau Yih, editors, *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 9594–9611, Online and Punta Cana, Dominican Republic, November 2021. Association for Computational Linguistics.
- 255 Xi Ye, Semih Yavuz, Kazuma Hashimoto, Yingbo Zhou, and Caiming Xiong. RNG-KBQA: Generation augmented iterative ranking for knowledge base question answering. In Smaranda Muresan, Preslav Nakov, and Aline Villavicencio, editors, *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 6032–6043, Dublin, Ireland, May 2022. Association for Computational Linguistics.
- 256 Tianle Li, Xueguang Ma, Alex Zhuang, Yu Gu, Yu Su, and Wenhui Chen. Few-shot in-context learning on knowledge base question answering. In Anna Rogers, Jordan Boyd-Graber, and Naoaki Okazaki, editors, *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 6966–6980, Toronto, Canada, July 2023. Association for Computational Linguistics.
- 257 Zhenyu Li, Sunqi Fan, Yu Gu, Xiuxing Li, Zhichao Duan, Bowen Dong, Ning Liu, and Jianyong Wang. Flexkbqa: A flexible llm-powered framework for few-shot knowledge base question answering. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 38, pages 18608–18616, 2024.
- 258 Haoran Luo, Haihong E, Zichen Tang, Shiyao Peng, Yikai Guo, Wentai Zhang, Chenghao Ma, Guanting Dong, Meina Song, Wei Lin, Yifan Zhu, and Anh Tuan Luu. ChatKBQA: A generate-then-retrieve framework for knowledge base question answering with fine-tuned large language models. In Lun-Wei Ku, Andre Martins, and Vivek Srikumar, editors, *Findings of the Association for Computational Linguistics: ACL 2024*, pages 2039–2056, Bangkok, Thailand, August 2024. Association for Computational Linguistics.
- 259 Zhiqiang Zhang, Liqiang Wen, and Wen Zhao. A gail fine-tuned llm enhanced framework for low-resource knowledge graph question answering. In *Proceedings of the 33rd ACM International Conference on Information and Knowledge Management, CIKM '24*, page 3300–3309, New York, NY, USA, 2024. Association for Computing Machinery.
- 260 Chang Zong, Yuchen Yan, Weiming Lu, Jian Shao, Yongfeng Huang, Heng Chang, and Yueting Zhuang. Triad: A framework leveraging a multi-role LLM-based agent to solve knowledge base question answering. In Yaser Al-Onaizan, Mohit Bansal, and Yun-Nung Chen, editors, *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, pages 1698–1710, Miami, Florida, USA, November 2024. Association for Computational Linguistics.
- 261 Apoorv Saxena, Aditay Tripathi, and Partha Talukdar. Improving multi-hop question answering over knowledge graphs using knowledge base embeddings. In Dan Jurafsky, Joyce Chai, Natalie Schluter, and Joel Tetreault, editors, *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 4498–4507, Online, July 2020. Association for Computational Linguistics.
- 262 Apoorv Saxena, Adrian Kochsiek, and Rainer Gemulla. Sequence-to-sequence knowledge graph completion and question answering. In Smaranda Muresan, Preslav Nakov, and Aline Villavicencio, editors, *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 2814–2828, Dublin, Ireland, May 2022. Association for Computational Linguistics.
- 263 Jing Zhang, Xiaokang Zhang, Jifan Yu, Jian Tang, Jie Tang, Cuiping Li, and Hong Chen. Subgraph retrieval enhanced model for multi-hop knowledge base question answering. In Smaranda Muresan, Preslav Nakov, and Aline Villavicencio, editors, *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 5773–5784, Dublin, Ireland, May 2022. Association for Computational Linguistics.
- 264 Jinhao Jiang, Kun Zhou, Xin Zhao, and Ji-Rong Wen. UniKGQA: Unified retrieval and reasoning for solving multi-hop question answering over knowledge graph. In *The Eleventh International Conference on Learning Representations*, 2023.
- 265 Donghan Yu, Sheng Zhang, Patrick Ng, Henghui Zhu, Alexander Hanbo Li, Jun Wang, Yiqun Hu, William Yang Wang, Zhiguo Wang, and Bing Xiang. DecAF: Joint decoding of answers and logical forms for question answering over knowledge bases. In *The Eleventh International Conference on Learning Representations*, 2023.
- 266 Linhao Luo, Yuan-Fang Li, Reza Haf, and Shirui Pan. Reasoning on graphs: Faithful and interpretable large language model reasoning. In *ICLR*, 2024.
- 267 Yu Gu, Xiang Deng, and Yu Su. Don't generate, discriminate: A proposal for grounding language models to real-world environments. In Anna Rogers, Jordan Boyd-Graber, and Naoaki Okazaki, editors, *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 4928–4949, Toronto, Canada, July 2023. Association for Computational Linguistics.
- 268 Yixin Ji, Kaixin Wu, Juntao Li, Wei Chen, Mingjie Zhong, Xu Jia, and Min Zhang. Retrieval and reasoning on KGs: Integrate knowledge graphs into large language models for complex question answering. In Yaser Al-Onaizan, Mohit Bansal, and Yun-Nung Chen, editors, *Findings of the Association for Computational Linguistics: EMNLP 2024*, pages 7598–7610, Miami, Florida, USA, November 2024. Association for Computational Linguistics.
- 269 Jinhao Jiang, Kun Zhou, Zican Dong, Keming Ye, Xin Zhao, and Ji-Rong Wen. StructGPT: A general framework for large language model to reason over structured data. In Houda Bouamor, Juan Pino, and Kalika Bali, editors, *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 9237–9251, Singapore, December 2023. Association for Computational Linguistics.
- 270 Jiashuo Sun, Chengjin Xu, Lumingyuan Tang, Saizhuo Wang, Chen Lin, Yeyun Gong, Lionel Ni, Heung-Yeung Shum, and Jian Guo. Think-on-graph: Deep and responsible reasoning of large language model on knowledge graph. In *The Twelfth International Conference on Learning Representations*, 2024.
- 271 Yike Wu, Yi Huang, Nan Hu, Yuncheng Hua, Guilin Qi, Jiaoyan Chen, and Jeff Z. Pan. CoTKR: Chain-of-thought enhanced knowledge rewriting for complex knowledge graph question answering. In Yaser Al-Onaizan, Mohit Bansal, and Yun-Nung Chen, editors, *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, pages 3501–3520, Miami, Florida, USA, November 2024. Association for Computational Linguistics.
- 272 Armin Toroghi, Willis Guo, Mohammad Mahdi Abdollah Pour, and Scott Sanner. Right for right reasons: Large language models for verifiable commonsense knowledge graph question answering. In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, pages 6601–6633. Association for Computational Linguistics, 2024.
- 273 Apoorv Saxena, Soumen Chakrabarti, and Partha Talukdar. Question answering over temporal knowledge graphs. In Chengqing Zong, Fei Xia, Wenjie Li, and Roberto Navigli, editors, *Proceedings of the 59th Annual Meeting of the As-*

- sociation for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers), pages 6663–6676, Online, August 2021. Association for Computational Linguistics.
- 274 Costas Mavromatis, Prasanna Lakkur Subramanyam, Vassilis N Ioannidis, Adesoji Adeshina, Phillip R Howard, Tetiana Grinberg, Nagib Hakim, and George Karypis. Tempqor: temporal question reasoning over knowledge graphs. In *Proceedings of the AAAI conference on artificial intelligence*, volume 36, pages 5825–5833, 2022.
- 275 Chao Shang, Guangtao Wang, Peng Qi, and Jing Huang. Improving time sensitivity for question answering over temporal knowledge graphs, 2022.
- 276 Yonghao Liu, Di Liang, Fang Fang, Sirui Wang, Wei Wu, and Rui Jiang. Time-aware multiway adaptive fusion network for temporal knowledge graph question answering. In *ICASSP 2023 - 2023 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 1–5, 2023.
- 277 Rikui Huang, Wei Wei, Xiaoye Qu, Wenfeng Xie, Xianling Mao, and Dangyang Chen. Joint multi-facts reasoning network for complex temporal question answering over knowledge graph, 2024.
- 278 Chenyang Du, Xiaoge Li, and Zhongyang Li. Semantic-enhanced reasoning question answering over temporal knowledge graphs. *Journal of Intelligent Information Systems*, pages 1–23, 2024.
- 279 Chao Xue, Di Liang, Pengfei Wang, and Jing Zhang. Question calibration and multi-hop modeling for temporal question answering. In *AAAI Conference on Artificial Intelligence*, 2024.
- 280 Zhen Jia, Soumajit Pramanik, Rishiraj Saha Roy, and Gerhard Weikum. Complex temporal question answering on knowledge graphs. In *Proceedings of the 30th ACM International Conference on Information & Knowledge Management, CIKM '21*, page 792–802, New York, NY, USA, 2021. Association for Computing Machinery.
- 281 Songlin Jiao, Zhenfang Zhu, Wenqing Wu, Zicheng Zuo, Jiangtao Qi, Wenling Wang, Guangyuan Zhang, and Peiyu Liu. An improving reasoning network for complex question answering over temporal knowledge graphs. *Applied Intelligence*, 53(7):8195–8208, July 2022.
- 282 Ziyang Chen, Xiang Zhao, Jinzhi Liao, Xinyi Li, and Evangelos Kanoulas. Temporal knowledge graph question answering via subgraph reasoning. *Knowledge-Based Systems*, 251:109134, 2022.
- 283 Aditya Sharma, Apoorv Saxena, Chitrang Gupta, Mehran Kazemi, Partha Talukdar, and Soumen Chakrabarti. TwiRGCN: Temporally weighted graph convolution for question answering over temporal knowledge graphs. In Andreas Vlachos and Isabelle Augenstein, editors, *Proceedings of the 17th Conference of the European Chapter of the Association for Computational Linguistics*, pages 2049–2060, Dubrovnik, Croatia, May 2023. Association for Computational Linguistics.
- 284 Yao Xiao, Guangyou Zhou, and Jin Liu. Modeling temporal-sensitive information for complex question answering over knowledge graphs. In Wei Lu, Shujian Huang, Yu Hong, and Xiabing Zhou, editors, *Natural Language Processing and Chinese Computing*, pages 418–430, Cham, 2022. Springer International Publishing.
- 285 Yonghao Liu, Di Liang, Mengyu Li, Fausto Giunchiglia, Ximing Li, Sirui Wang, Wei Wu, Lan Huang, Xiaoyue Feng, and Renchu Guan. Local and global: Temporal question answering via information fusion. In Edith Elkind, editor, *Proceedings of the Thirty-Second International Joint Conference on Artificial Intelligence, IJCAI-23*, pages 5141–5149. International Joint Conferences on Artificial Intelligence Organization, 8 2023. Main Track.
- 286 Zhiyuan Zha, Pengnian Qi, Xigang Bao, Mengyuan Tian, and Biao Qin. M3tqa: Multi-view, multi-hop and multi-stage reasoning for temporal question answering. In *ICASSP 2024 - 2024 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 10086–10090, 2024.
- 287 Zhuo Chen, Zhao Zhang, Zixuan Li, Fei Wang, Yutao Zeng, Xiaolong Jin, and Yongjun Xu. Self-improvement programming for temporal knowledge graph question answering. In Nicoletta Calzolari, Min-Yen Kan, Veronique Hoste, Alessandro Lenci, Sakriani Sakti, and Nianwen Xue, editors, *Proceedings of the 2024 Joint International Conference on Computational Linguistics, Language Resources and Evaluation (LREC-COLING 2024)*, pages 14579–14594, Torino, Italia, May 2024. ELRA and ICCL.
- 288 Ziyang Chen, Dongfang Li, Xiang Zhao, Baotian Hu, and Min Zhang. Temporal knowledge question answering via abstract reasoning induction. In Lun-Wei Ku, Andre Martins, and Vivek Srikumar, editors, *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 4872–4889, Bangkok, Thailand, August 2024. Association for Computational Linguistics.
- 289 Xinying Qian, Ying Zhang, Yu Zhao, Baohang Zhou, Xuhui Sui, Li Zhang, and Kehui Song. TimeR⁴: Time-aware retrieval-augmented large language models for temporal knowledge graph question answering. In Yaser Al-Onaizan, Mohit Bansal, and Yun-Nung Chen, editors, *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, pages 6942–6952, Miami, Florida, USA, November 2024. Association for Computational Linguistics.
- 290 Jonathan Herzig, Pawel Krzysztof Nowak, Thomas Müller, Francesco Piccinno, and Julian Eisenschlos. TaPas: Weakly supervised table parsing via pre-training. In Dan Jurafsky, Joyce Chai, Natalie Schluter, and Joel Tetreault, editors, *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 4320–4333, Online, July 2020. Association for Computational Linguistics.
- 291 Tao Yu, Chien-Sheng Wu, Xi Victoria Lin, hailin wang, Yi Chern Tan, Xinyi Yang, Dragomir Radev, richard socher, and Caiming Xiong. Gra{pp}a: Grammar-augmented pre-training for table semantic parsing. In *International Conference on Learning Representations*, 2021.
- 292 Rami Aly and Andreas Vlachos. Tabver: Tabular fact verification with natural logic, 2024.
- 293 Wanjun Zhong, Junjie Huang, Qian Liu, Ming Zhou, Jiahai Wang, Jian Yin, and Nan Duan. Reasoning over hybrid chain for table-and-text open domain question answering. In Lud De Raedt, editor, *Proceedings of the Thirty-First International Joint Conference on Artificial Intelligence, IJCAI-22*, pages 4531–4537. International Joint Conferences on Artificial Intelligence Organization, 7 2022. Main Track.
- 294 Haowei Zhang, Shengyuan Si, Yilun Zhao, Lujing Xie, Zhijian Xu, Lyuhao Chen, Linyong Nan, Pengcheng Wang, Xiangru Tang, and Arman Cohan. OpenT2T: An open-source toolkit for table-to-text generation. In Delia Irazu Hernandez Farias, Tom Hope, and Manling Li, editors, *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing: System Demonstrations*, pages 259–269, Miami, Florida, USA, November 2024. Association for Computational Linguistics.
- 295 Peng Li, Yeye He, Dror Yashar, Weiwei Cui, Song Ge, Haidong Zhang, Danielle Rifinski Fainman, Dongmei Zhang, and Surajit Chaudhuri. Table-gpt: Table fine-tuned gpt for diverse table tasks. *Proc. ACM Manag. Data*, 2(3), May 2024.
- 296 Victor Zhong, Caiming Xiong, and Richard Socher. Seq2sql: Generating structured queries from natural language using reinforcement learning, 2017.
- 297 Zirui Wu and Yansong Feng. Protrix: Building models for planning and reasoning over tables with sentence context, 2024.
- 298 Tianshu Zhang, Xiang Yue, Yifei Li, and Huan Sun. TableLlama: Towards open large generalist models for tables. In Kevin Duh, Helena Gomez, and Steven Bethard, editors, *Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)*, pages 6024–6044, Mexico City, Mexico, June 2024. Association for Computational Linguistics.
- 299 Zihan Liu, Wei Ping, Rajarshi Roy, Peng Xu, Chankyu Lee, Mohammad Shoeybi, and Bryan Catanzaro. ChatQA: Surpassing GPT-4 on conversational QA and RAG. In *The Thirty-eighth Annual Conference on Neural Information Processing Systems*, 2024.
- 300 Vaishali Pal, Evangelos Kanoulas, Andrew Yates, and Maarten de Rijke. Table question answering for low-resourced Indic

- languages. In Yaser Al-Onaizan, Mohit Bansal, and Yun-Nung Chen, editors, *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, pages 75–92, Miami, Florida, USA, November 2024. Association for Computational Linguistics.
- 301 Fengbin Zhu, Wenqiang Lei, Youcheng Huang, Chao Wang, Shuo Zhang, Jiancheng Lv, Fuli Feng, and Tat-Seng Chua. TAT-QA: A question answering benchmark on a hybrid of tabular and textual content in finance. In Chengqing Zong, Fei Xia, Wenjie Li, and Roberto Navigli, editors, *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pages 3277–3287, Online, August 2021. Association for Computational Linguistics.
- 302 Fangyu Lei, Xiang Li, Yifan Wei, Shizhu He, Yiming Huang, Jun Zhao, and Kang Liu. S3HQA: A three-stage approach for multi-hop text-table hybrid question answering. In Anna Rogers, Jordan Boyd-Graber, and Naoaki Okazaki, editors, *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, pages 1731–1740, Toronto, Canada, July 2023. Association for Computational Linguistics.
- 303 Tongxu Luo, Fangyu Lei, Jiahe Lei, Weihao Liu, Shihu He, Jun Zhao, and Kang Liu. Hrot: Hybrid prompt strategy and retrieval of thought for table-text hybrid question answering, 2023.
- 304 Weihao Liu, Fangyu Lei, Tongxu Luo, Jiahe Lei, Shizhu He, Jun Zhao, and Kang Liu. Mmhqa-icl: Multimodal in-context learning for hybrid question answering over text, tables and images, 2023.
- 305 Kiran Purohit, Venkatesh V, Raghuram Devalla, Krishna Mohan Yerragorla, Sourangshu Bhattacharya, and Avishek Anand. Explora: Efficient exemplar subset selection for complex reasoning, 2024.
- 306 Xianjie Wu, Jian Yang, Linzheng Chai, Ge Zhang, Jiaheng Liu, Xinrun Du, Di Liang, Daixin Shu, Xianfu Cheng, Tianzhen Sun, Guanglin Niu, Tongliang Li, and Zhoujun Li. Tablebench: A comprehensive and complex benchmark for table question answering, 2024.
- 307 Mengkang Hu, Haoyu Dong, Ping Luo, Shi Han, and Dongmei Zhang. KET-QA: A dataset for knowledge enhanced table question answering. In Nicoletta Calzolari, Min-Yen Kan, Veronique Hoste, Alessandro Lenci, Sakriani Sakti, and Nianwen Xue, editors, *Proceedings of the 2024 Joint International Conference on Computational Linguistics, Language Resources and Evaluation (LREC-COLING 2024)*, pages 9705–9719, Torino, Italia, May 2024. ELRA and ICCL.
- 308 Yasuhito Ohsugi Itsumi Saito Kyosuke Nishida and Hisako Asano Junji Tomita. A simple but effective method to incorporate multi-turn context with bert for conversational machine comprehension. *ACL 2019*, page 11, 2019.
- 309 Zhuosheng Zhang, Yuwei Wu, Junru Zhou, Sufeng Duan, Hai Zhao, and Rui Wang. Sg-net: Syntax-guided machine reading comprehension. In *Proceedings of the AAAI conference on artificial intelligence*, volume 34, pages 9636–9643, 2020.
- 310 Elad Segal, Avia Efrat, Mor Shoham, Amir Globerson, and Jonathan Berant. A simple and effective model for answering multi-span questions. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 3074–3080, 2020.
- 311 Minghao Hu, Yuxing Peng, Zhen Huang, and Dongsheng Li. A multi-type multi-span network for reading comprehension that requires discrete reasoning. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 1596–1606, 2019.
- 312 Dan Su, Xiaoguang Li, Jindi Zhang, Lifeng Shang, Xin Jiang, Qun Liu, and Pascale Fung. Read before generate! faithful long form question answering with machine reading. In *60th Annual Meeting of the Association-for-Computational-Linguistics (ACL)*, pages 744–756. Association for Computational Linguistics, 2022.
- 313 Bin Bi, Chen Wu, Ming Yan, Wei Wang, Jiangnan Xia, and Chenliang Li. Incorporating external knowledge into machine reading for generative question answering. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 2521–2530, 2019.
- 314 Chenliang Li, Bin Bi, Ming Yan, Wei Wang, and Songfang Huang. Addressing semantic drift in generative question answering with auxiliary extraction. In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 2: Short Papers)*, pages 942–947, 2021.
- 315 Peng Xu, Davis Liang, Zhiheng Huang, and Bing Xiang. Attention-guided generative models for extractive question answering. *arXiv preprint arXiv:2110.06393*, 2021.
- 316 Lin Ai, Zheng Hui, Zizhou Liu, and Julia Hirschberg. Enhancing pre-trained generative language models with question attended span extraction on machine reading comprehension. In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, pages 10046–10063. Association for Computational Linguistics, 2024.
- 317 Junjie Yang, Zhuosheng Zhang, and Hai Zhao. Multi-span style extraction for generative reading comprehension. *arXiv preprint arXiv:2009.07382*, 2020.
- 318 Zhuosheng Zhang, Yiqing Zhang, Hai Zhao, Xi Zhou, and Xiang Zhou. Composing answer from multi-spans for reading comprehension. *arXiv preprint arXiv:2009.06141*, 2020.
- 319 Zhuosheng Zhang, Yiqing Zhang, and Hai Zhao. Syntax-aware multi-spans generation for reading comprehension. *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, 30:260–268, 2021.
- 320 Ming Ding, Chang Zhou, Qibin Chen, Hongxia Yang, and Jie Tang. Cognitive graph for multi-hop reading comprehension at scale. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 2694–2703, 2019.
- 321 Xiaohui Li, Yuezhong Liu, Shenggen Ju, and Zhengwen Xie. Dynamic reasoning network for multi-hop question answering. In *Natural Language Processing and Chinese Computing: 9th CCF International Conference, NLPCC 2020, Zhengzhou, China, October 14–18, 2020, Proceedings, Part I 9*, pages 29–40. Springer, 2020.
- 322 Alex Long, Joel Mason, Alan Blair, and Wei Wang. Multi-hop reading comprehension via deep reinforcement learning based document traversal. *arXiv preprint arXiv:1905.09438*, 2019.
- 323 Jifan Chen, Shih-ting Lin, and Greg Durrett. Multi-hop question answering via reasoning chains. *arXiv preprint arXiv:1910.02610*, 2019.
- 324 Jiuyang Tang, Shengze Hu, Ziyang Chen, Hao Xu, and Zhen Tan. Incorporating phrases in latent query reformulation for multi-hop question answering. *Mathematics*, 10(4):646, 2022.
- 325 Ming Tu, Guangtao Wang, Jing Huang, Yun Tang, Xiaodong He, and Bowen Zhou. Multi-hop reading comprehension across multiple documents by reasoning over heterogeneous graphs. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*. Association for Computational Linguistics, 2019.
- 326 Yuwei Fang, Siqi Sun, Zhe Gan, Rohit Pillai, Shuohang Wang, and Jingjing Liu. Hierarchical graph network for multi-hop question answering. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 8823–8838, 2020.
- 327 Azade Mohammadi, Reza Ramezani, and Ahmad Baraani. Topic-aware multi-hop machine reading comprehension using weighted graphs. *Expert Systems with Applications*, 224:119873, 2023.
- 328 Nan Shao, Yiming Cui, Ting Liu, Shijin Wang, and Guoping Hu. Is graph structure necessary for multi-hop question answering? In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 7187–7192, 2020.
- 329 Sewon Min, Victor Zhong, Luke Zettlemoyer, and Hannaneh Hajishirzi. Multi-hop reading comprehension through question

- decomposition and rescoring. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 6097–6109, 2019.
- 330 Dirk Groeneveld, Tushar Khot, Ashish Sabharwal, et al. A simple yet strong pipeline for hotpotqa. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 8839–8845, 2020.
- 331 Bohong Wu, Zhuosheng Zhang, and Hai Zhao. Graph-free multi-hop reading comprehension: A select-to-guide strategy. *arXiv preprint arXiv:2107.11823*, 2021.
- 332 Yin Zhangyue, Wang Yuxin, Hu Xiannian, Wu Yiguang, Yan Hang, Zhang Xinyu, Cao Zhao, Huang Xuanjing, and Qiu Xipeng. Rethinking label smoothing on multi-hop question answering. In *Proceedings of the 22nd Chinese National Conference on Computational Linguistics*, pages 611–623, 2023.
- 333 Xin-Yi Li, Wei-Jun Lei, and Yu-Bin Yang. From easy to hard: Two-stage selector and reader for multi-hop question answering. In *ICASSP 2023-2023 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 1–5. IEEE, 2023.
- 334 Yung-Sung Chuang, Wei Fang, Shang-Wen Li, Wen-tau Yih, and James Glass. Expand, rerank, and retrieve: Query reranking for open-domain question answering. In *Findings of the Association for Computational Linguistics: ACL 2023*, pages 12131–12147, 2023.
- 335 Minsang Kim, Cheoneum Park, and Seung Baek. Qpaug: Question and passage augmentation for open-domain question answering of llms. In *Findings of the Association for Computational Linguistics: EMNLP 2024*, pages 9024–9042, 2024.
- 336 Chunlei Xin, Yaojie Lu, Hongyu Lin, Shuheng Zhou, Huijia Zhu, Weiqiang Wang, Zhongyi Liu, Xianpei Han, and Le Sun. Chain-of-rewrite: Aligning question and documents for open-domain question answering. In *Findings of the Association for Computational Linguistics: EMNLP 2024*, pages 1884–1896, 2024.
- 337 Jamshid Mozafari, Abdelrahman Abdallah, Bhawna Piryani, and Adam Jatowt. Exploring hint generation approaches for open-domain question answering. In *Findings of the Association for Computational Linguistics: EMNLP 2024*, pages 9327–9352, 2024.
- 338 Wenhao Yu, Dan Iter, Shuohang Wang, Yichong Xu, Mingxuan Ju, S Sanyal, Chenguang Zhu, Michael Zeng, and Meng Jiang. Generate rather than retrieve: Large language models are strong context generators. In *International Conference on Learning Representations*, 2023.
- 339 Junlong Li, Jinyuan Wang, Zhuosheng Zhang, and Hai Zhao. Self-prompting large language models for zero-shot open-domain QA. In *Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)*, pages 296–310, 2024.
- 340 Giacomo Frisoni, Alessio Cocchieri, Alex Presepi, Gianluca Moro, and Zaiqiao Meng. To generate or to retrieve? on the effectiveness of artificial contexts for medical open-domain question answering. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 9878–9919, August 2024.
- 341 Sukmin Cho, Jeongyeon Seo, Soyeong Jeong, and Jong C Park. Improving zero-shot reader by reducing distractions from irrelevant documents in open-domain question answering. In *Findings of the Association for Computational Linguistics: EMNLP 2023*, pages 3145–3157, 2023.
- 342 Yuhao Wang, Ruiyang Ren, Junyi Li, Xin Zhao, Jing Liu, and Ji-Rong Wen. REAR: A relevance-aware retrieval-augmented framework for open-domain question answering. In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, pages 5613–5626, 2024.
- 343 Kiseung Kim and Jay-Yoon Lee. RE-RAG: Improving open-domain QA performance and interpretability with relevance estimator in retrieval-augmented generation. In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, pages 22149–22161, 2024.
- 344 Yufei Huang, Xu Han, and Maosong Sun. Fastfid: Improve inference efficiency of open domain question answering via sentence selection. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 6262–6276, 2024.
- 345 Zihan Zhang, Meng Fang, and Ling Chen. RetrievalQA: Assessing adaptive retrieval-augmented generation for short-form open-domain question answering. In *Findings of the Association for Computational Linguistics: ACL 2024*, pages 6963–6975, August 2024.
- 346 Yu Wang, Nedim Lipka, Ryan A Rossi, Alexa Siu, Ruiyi Zhang, and Tyler Derr. Knowledge graph prompting for multi-document question answering. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 38, pages 19206–19214, 2024.
- 347 Zukang Yang and Zixuan Zhu. Curiousllm: Elevating multi-document qa with reasoning-infused knowledge graph prompting. *arXiv preprint arXiv:2404.09077*, 2024.
- 348 Barlas Oguz, Xilun Chen, Vladimir Karpukhin, Stan Peshterliev, Dmytro Okhonko, Michael Schlichtkrull, Sonal Gupta, Yashar Mehdad, and Scott Yih. UniK-QA: Unified representations of structured and unstructured knowledge for open-domain question answering. In *Findings of the Association for Computational Linguistics: NAACL 2022*, pages 1535–1546, 2022.
- 349 Kaixin Ma, Hao Cheng, Xiaodong Liu, Eric Nyberg, and Jianfeng Gao. Open domain question answering with a unified knowledge interface. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1605–1620, 2022.
- 350 Philipp Christmann, Rishiraj Saha Roy, and Gerhard Weikum. Conversational question answering on heterogeneous sources. In *Proceedings of the 45th International ACM SIGIR Conference on Research and Development in Information Retrieval*, pages 144–154, 2022.
- 351 Wen Zhang, Long Jin, Yushan Zhu, Jiaoyan Chen, Zhiwei Huang, Junjie Wang, Yin Hua, Lei Liang, and Huajun Chen. Trustuqa: A trustful framework for unified structured data question answering. *arXiv preprint arXiv:2406.18916*, 2024.
- 352 Philipp Christmann, Rishiraj Saha Roy, and Gerhard Weikum. Explainable conversational question answering over heterogeneous sources via iterative graph neural networks. In *Proceedings of the 46th International ACM SIGIR Conference on Research and Development in Information Retrieval*, pages 643–653, 2023.
- 353 Heidi Zhang, Sina Semnani, Farhad Ghassemi, Jialiang Xu, Shicheng Liu, and Monica Lam. SPAGHETTI: Open-domain question answering from heterogeneous data sources with retrieval and semantic parsing. In *Findings of the Association for Computational Linguistics: ACL 2024*, pages 1663–1678, 2024.
- 354 Hugo Touvron, Louis Martin, Kevin R. Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shrutu Bhosale, Daniel M. Bikel, Lukas Blecher, Cristian Cantón Ferrer, Moya Chen, Guillem Cucurull, David Esiobu, Jude Fernandes, Jeremy Fu, Wenyin Fu, Brian Fuller, Cynthia Gao, Vedanuj Goswami, Naman Goyal, Anthony S. Hartshorn, Saghar Hosseini, Rui Hou, Hakan Inan, Marcin Kardas, Viktor Kerkez, Madian Khabsa, Isabel M. Kloumann, A. V. Korenev, Punit Singh Koura, Marie-Anne Lachaux, Thibaut Lavril, Jenya Lee, Diana Liskovich, Yinghai Lu, Yuning Mao, Xavier Martinet, Todor Mihaylov, Pushkar Mishra, Igor Molybog, Yixin Nie, Andrew Poulton, Jeremy Reizenstein, Rashi Rungta, Kalyan Saladi, Alan Schelten, Ruan Silva, Eric Michael Smith, R. Subramanian, Xia Tan, Binh Tang, Ross Taylor, Adina Williams, Jian Xiang Kuan, Puxin Xu, Zhengxu Yan, Iliyan Zarov, Yuchen Zhang, Angela Fan, Melanie Kambadur, Sharan Narang, Aurelien Rodriguez, Robert Stojnic, Sergey Edunov, and Thomas Scialom. Llama 2: Open foundation and fine-tuned chat models. *ArXiv*, abs/2307.09288, 2023.
- 355 Gautier Izacard and Edouard Grave. Leveraging passage retrieval with generative models for open domain question answer-

- ing. In Paola Merlo, Jorg Tiedemann, and Reut Tsarfaty, editors, *Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics: Main Volume*, pages 874–880, Online, April 2021. Association for Computational Linguistics.
- 356 Nils Reimers and Iryna Gurevych. Sentence-BERT: Sentence embeddings using Siamese BERT-networks. In Kentaro Inui, Jing Jiang, Vincent Ng, and Xiaojun Wan, editors, *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 3982–3992, Hong Kong, China, November 2019. Association for Computational Linguistics.
- 357 Yifu Gao, Linbo Qiao, Zhigang Kan, Zhihua Wen, Yongquan He, and Dongsheng Li. Two-stage generative question answering on temporal knowledge graph using large language models. In Lun-Wei Ku, Andre Martins, and Vivek Srikumar, editors, *Findings of the Association for Computational Linguistics: ACL 2024*, pages 6719–6734, Bangkok, Thailand, August 2024. Association for Computational Linguistics.
- 358 Qingqing Cai and Alexander Yates. Large-scale semantic parsing via schema matching and lexicon extension. In Hinrich Schuetze, Pascale Fung, and Massimo Poesio, editors, *Proceedings of the 51st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 423–433, Sofia, Bulgaria, August 2013. Association for Computational Linguistics.
- 359 Jonathan Berant, Andrew Chou, Roy Frostig, and Percy Liang. Semantic parsing on Freebase from question-answer pairs. In David Yarowsky, Timothy Baldwin, Anna Korhonen, Karen Livescu, and Steven Bethard, editors, *Proceedings of the 2013 Conference on Empirical Methods in Natural Language Processing*, pages 1533–1544, Seattle, Washington, USA, October 2013. Association for Computational Linguistics.
- 360 Wen-tau Yih, Matthew Richardson, Chris Meek, Ming-Wei Chang, and Jina Suh. The value of semantic parse labeling for knowledge base question answering. In Katrin Erk and Noah A. Smith, editors, *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, pages 201–206, Berlin, Germany, August 2016. Association for Computational Linguistics.
- 361 Junwei Bao, Nan Duan, Zhao Yan, Ming Zhou, and Tiejun Zhao. Constraint-based question answering with knowledge graph. In Yuji Matsumoto and Rashmi Prasad, editors, *Proceedings of COLING 2016, the 26th International Conference on Computational Linguistics: Technical Papers*, pages 2503–2514, Osaka, Japan, December 2016. The COLING 2016 Organizing Committee.
- 362 Yuyu Zhang, Hanjun Dai, Zornitsa Kozareva, Alexander Smola, and Le Song. Variational reasoning for question answering with knowledge graph. In *Proceedings of the AAAI conference on artificial intelligence*, volume 32, 2018.
- 363 Ricardo Usbeck, Xiongliang Yan, Aleksandr Perevalov, Longquan Jiang, Julius Schulz, Angelie Kraft, Cedric Möller, Junbo Huang, Jan Reineke, Axel-Cyrille Ngonga Ngomo, Muhammad Saleem, and Andreas Both. Qald-10 – the 10th challenge on question answering over linked data. *Semantic Web*, 2023.
- 364 Priyansh Trivedi, Gaurav Maheshwari, Mohnish Dubej, and Jens Lehmann. Lc-quad: A corpus for complex question answering over knowledge graphs. In *The Semantic Web – ISWC 2017: 16th International Semantic Web Conference, Vienna, Austria, October 21–25, 2017, Proceedings, Part II*, page 210–218, Berlin, Heidelberg, 2017. Springer-Verlag.
- 365 Mohnish Dubej, Debayan Banerjee, Abdelrahman Abdelkawi, and Jens Lehmann. Lc-quad 2.0: A large dataset for complex question answering over wikidata and dbpedia. In Chiara Ghidini, Olaf Hartig, Maria Maleshkova, Vojtěch Svátek, Isabel Cruz, Aidan Hogan, Jie Song, Maxime Lefrançois, and Fabien Gandon, editors, *The Semantic Web – ISWC 2019*, pages 69–78, Cham, 2019. Springer International Publishing.
- 366 Zhen Jia, Abdalghani Abujabal, Rishiraj Saha Roy, Jannik Strötgen, and Gerhard Weikum. Tempquestions: A benchmark for temporal question answering. In *Companion Proceedings of the The Web Conference 2018*, WWW '18, page 1057–1062, Republic and Canton of Geneva, CHE, 2018. International World Wide Web Conferences Steering Committee.
- 367 Ziyang Chen, Jinzhi Liao, and Xiang Zhao. Multi-granularity temporal question answering over knowledge graphs. In Anna Rogers, Jordan Boyd-Graber, and Naoaki Okazaki, editors, *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 11378–11392, Toronto, Canada, July 2023. Association for Computational Linguistics.
- 368 Thomas Mueller, Francesco Piccinno, Peter Shaw, Massimo Nicosia, and Yasemin Altun. Answering conversational questions on structured data without logical forms. In Kentaro Inui, Jing Jiang, Vincent Ng, and Xiaojun Wan, editors, *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 5902–5910, Hong Kong, China, November 2019. Association for Computational Linguistics.
- 369 Kushal Raj Bhandari, Sixue Xing, Soham Dan, and Jianxi Gao. On the robustness of language models for tabular question answering, 2024.
- 370 Sam Wiseman, Stuart Shieber, and Alexander Rush. Challenges in data-to-document generation. In *Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing*, pages 2253–2263, 2017.
- 371 Lya Hullyyyatus Suadaa, Hidetaka Kamigaito, Kotaro Funakoshi, Manabu Okumura, and Hiroya Takamura. Towards table-to-text generation with numerical reasoning. In Chengqing Zong, Fei Xia, Wenjie Li, and Roberto Navigli, editors, *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pages 1451–1465, Online, August 2021. Association for Computational Linguistics.
- 372 Nafise Sadat Moosavi, Andreas Rücklé, Dan Roth, and Iryna Gurevych. Scigen: a dataset for reasoning-aware text generation from scientific tables. In *Thirty-fifth Conference on Neural Information Processing Systems Datasets and Benchmarks Track (Round 2)*, 2021.
- 373 Wenhui Chen, Jianshu Chen, Yu Su, Zhiyu Chen, and William Yang Wang. Logical natural language generation from open-domain tables. In Dan Jurafsky, Joyce Chai, Natalie Schluter, and Joel Tetreault, editors, *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 7929–7942, Online, July 2020. Association for Computational Linguistics.
- 374 Qixiang Zhou, Tong Li, and Yunduo Wang. Assisting in requirements goal modeling: a hybrid approach based on machine learning and logical reasoning. In *Proceedings of the 25th International Conference on Model Driven Engineering Languages and Systems, MODELS '22*, page 199–209, New York, NY, USA, 2022. Association for Computing Machinery.
- 375 Frank van Harmelen and Annette ten Teije. A boxology of design patterns for hybrid learning and reasoning systems. *Journal of Web Engineering*, 18(1-3):97–123, 2019.
- 376 Zhoujun Cheng, Haoyu Dong, Zhiruo Wang, Ran Jia, Jiaqi Guo, Yan Gao, Shi Han, Jian-Guang Lou, and Dongmei Zhang. HiTab: A hierarchical table dataset for question answering and natural language generation. In Smaranda Muresan, Preslav Nakov, and Aline Villavicencio, editors, *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1094–1110, Dublin, Ireland, May 2022. Association for Computational Linguistics.
- 377 Linyong Nan, Chiachun Hsieh, Ziming Mao, Xi Victoria Lin, Neha Verma, Rui Zhang, Wojciech Kryściński, Hailey Schoelkopf, Riley Kong, Xiangru Tang, Mutethia Mutuma, Ben Rosand, Isabel Trindade, Renusree Bandaru, Jacob Cunningham, Caiming Xiong, Dragomir Radev, and Dragomir Radev. Fetaqa: Free-form table question answering. *Transactions of the Association for Computational Linguistics*, 10:35–49, 01 2022.

- 378 Rami Aly, Zhijiang Guo, Michael Sejr Schlichtkrull, James Thorne, Andreas Vlachos, Christos Christodoulopoulos, Oana Cocarascu, and Arpit Mittal. The fact extraction and VERification over unstructured and structured information (FEVEROUS) shared task. In Rami Aly, Christos Christodoulopoulos, Oana Cocarascu, Zhijiang Guo, Arpit Mittal, Michael Schlichtkrull, James Thorne, and Andreas Vlachos, editors, *Proceedings of the Fourth Workshop on Fact Extraction and VERification (FEVER)*, pages 1–13, Dominican Republic, November 2021. Association for Computational Linguistics.
- 379 Danqi Chen, Jason Bolton, and Christopher D Manning. A thorough examination of the cnn/daily mail reading comprehension task. In *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 2358–2367, 2016.
- 380 Mrinmaya Sachan and Eric Xing. Machine comprehension using rich semantic representations. In *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, pages 486–492, 2016.
- 381 Moxin Li, Wenjie Wang, Fuli Feng, Hanwang Zhang, Qifan Wang, and Tat-Seng Chua. Hypothetical training for robust machine reading comprehension of tabular context. In *Findings of the Association for Computational Linguistics: ACL 2023*, pages 1220–1236, 2023.
- 382 Felix Hill, Antoine Bordes, Sumit Chopra, and Jason Weston. The goldilocks principle: Reading children’s books with explicit memory representations. In *4th International Conference on Learning Representations, ICLR 2016*, 2016.
- 383 Karl Moritz Hermann, Tomas Kocisky, Edward Grefenstette, Lasse Espeholt, Will Kay, Mustafa Suleyman, and Phil Blunsom. Teaching machines to read and comprehend. *Advances in neural information processing systems*, 28, 2015.
- 384 Pranav Rajpurkar, Jian Zhang, Konstantin Lopyrev, and Percy Liang. SQuAD: 100,000+ questions for machine comprehension of text. In Jian Su, Kevin Duh, and Xavier Carreras, editors, *Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing*, pages 2383–2392, Austin, Texas, November 2016. Association for Computational Linguistics.
- 385 Mandar Joshi, Eunsol Choi, Daniel S Weld, and Luke Zettlemoyer. Triviaqa: A large scale distantly supervised challenge dataset for reading comprehension. In *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1601–1611, 2017.
- 386 Guokun Lai, Qizhe Xie, Hanxiao Liu, Yiming Yang, and Eduard Hovy. Race: Large-scale reading comprehension dataset from examinations. In *Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing*, pages 785–794, 2017.
- 387 Payal Bajaj, Daniel Campos, Nick Craswell, Li Deng, Jianfeng Gao, Xiaodong Liu, Rangan Majumder, Andrew McNamara, Bhaskar Mitra, Tri Nguyen, et al. Ms marco: A human generated machine reading comprehension dataset. *arXiv preprint arXiv:1611.09268*, 2016.
- 388 Lynette Hirschman, Marc Light, Eric Breck, and John D Burger. Deep read: A reading comprehension system. In *Proceedings of the 37th annual meeting of the Association for Computational Linguistics*, pages 325–332, 1999.
- 389 Ellen Riloff and Michael Thelen. A rule-based question answering system for reading comprehension tests. In *ANLP-NAACL 2000 workshop: reading comprehension tests as evaluation for computer-based language understanding systems*, 2000.
- 390 Eugene Charniak, Yasemin Altun, Rodrigo de Salvo Braz, Benjamin Garrett, Margaret Kosmala, Tomer Moscovich, Lixin Pang, Changhee Pyo, Ye Sun, Wei Wy, et al. Reading comprehension programs in a statistical-language-processing class. In *ANLP-NAACL 2000 Workshop: Reading Comprehension Tests as Evaluation for Computer-Based Language Understanding Systems*, 2000.
- 391 Razieh Baradaran and Hossein Amirkhani. Ensemble learning-based approach for improving generalization capability of machine reading comprehension systems. *Neurocomputing*, 466:229–242, 2021.
- 392 Minghao Hu, Furu Wei, Yuxing Peng, Zhen Huang, Nan Yang, and Dongsheng Li. Read+ verify: Machine reading comprehension with unanswerable questions. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 33, pages 6529–6537, 2019.
- 393 Rudolf Kadlec, Martin Schmid, Ondřej Bajgar, and Jan Kleindienst. Text understanding with the attention sum reader network. In *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 908–918, 2016.
- 394 Bhuwan Dhingra, Hanxiao Liu, Zhilin Yang, William Cohen, and Ruslan Salakhutdinov. Gated-attention readers for text comprehension. In *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1832–1846, 2017.
- 395 Yiming Cui, Zhipeng Chen, Si Wei, Shijin Wang, Ting Liu, and Guoping Hu. Attention-over-attention neural networks for reading comprehension. In *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 593–602, 2017.
- 396 Min Joon Seo, Aniruddha Kembhavi, Ali Farhadi, and Hannaneh Hajishirzi. Bidirectional attention flow for machine comprehension. In *5th International Conference on Learning Representations, ICLR 2017, Toulon, France, April 24–26, 2017, Conference Track Proceedings*. OpenReview.net, 2017.
- 397 Zhilin Yang, Peng Qi, Saizheng Zhang, Yoshua Bengio, William Cohen, Ruslan Salakhutdinov, and Christopher D Manning. Hotpotqa: A dataset for diverse, explainable multi-hop question answering. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*. Association for Computational Linguistics, 2018.
- 398 Nicola De Cao, Wilker Aziz, and Ivan Titov. Question answering by reasoning across documents with graph convolutional networks. In *Proceedings of NAACL-HLT*, pages 2306–2317, 2019.
- 399 Yu Cao, Meng Fang, and Dacheng Tao. Bag: Bi-directional attention entity graph convolutional network for multi-hop reasoning question answering. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 357–362, 2019.
- 400 Lin Qiu, Yunxuan Xiao, Yanru Qu, Hao Zhou, Lei Li, Weinan Zhang, and Yong Yu. Dynamically fused graph network for multi-hop reasoning. In *Proceedings of the 57th annual meeting of the association for computational linguistics*, pages 6140–6150, 2019.
- 401 Xing Cao, Yun Liu, Bo Hu, and Yu Zhang. Dual-channel reasoning model for complex question answering. *Complexity*, 2021(1):7367181, 2021.
- 402 Ming Tu, Kevin Huang, Guangtao Wang, Jing Huang, Xiaodong He, and Bowen Zhou. Select, answer and explain: Interpretable multi-hop reading comprehension over multiple documents. In *Proceedings of the AAAI conference on artificial intelligence*, volume 34, pages 9073–9080, 2020.
- 403 Jonathan St BT Evans. Heuristic and analytic processes in reasoning. *British Journal of Psychology*, 75(4):451–468, 1984.
- 404 Jonathan St BT Evans. In two minds: dual-process accounts of reasoning. *Trends in cognitive sciences*, 7(10):454–459, 2003.
- 405 Jonathan St BT Evans. Dual-processing accounts of reasoning, judgment, and social cognition. *Annu. Rev. Psychol.*, 59(1):255–278, 2008.
- 406 Steven A Sloman. The empirical case for two systems of reasoning. *Psychological bulletin*, 119(1):3, 1996.
- 407 Yoshua Bengio, Jérôme Louradour, Ronan Collobert, and Jason Weston. Curriculum learning. In *Proceedings of the 26th annual international conference on machine learning*, pages 41–48, 2009.
- 408 Danqi Chen, Adam Fisch, Jason Weston, and Antoine Bordes. Reading wikipedia to answer open-domain questions. In

- Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, Jan 2017.
- 409 Gerard Salton, Edward A Fox, and Harry Wu. Extended boolean information retrieval. *Communications of the ACM*, 26(11):1022–1036, 1983.
- 410 Stephen Robertson, Hugo Zaragoza, et al. The probabilistic relevance framework: Bm25 and beyond. *Foundations and Trends® in Information Retrieval*, 3(4):333–389, 2009.
- 411 Thibault Formal, Benjamin Piwowarski, and Stéphane Clinchant. Splade: Sparse lexical and expansion model for first stage ranking. In *Proceedings of the 44th International ACM SIGIR Conference on Research and Development in Information Retrieval*, pages 2288–2292, 2021.
- 412 Jimmy Lin and Xueguang Ma. A few brief notes on deepimpact, coil, and a conceptual framework for information retrieval techniques. *arXiv preprint arXiv:2106.14807*, 2021.
- 413 Vladimir Karpukhin, Barlas Oğuz, Sewon Min, Patrick Lewis, Ledell Wu, Sergey Edunov, Danqi Chen, and Wen Tau Yih. Dense passage retrieval for open-domain question answering. In *2020 Conference on Empirical Methods in Natural Language Processing, EMNLP 2020*, pages 6769–6781. Association for Computational Linguistics (ACL), 2020.
- 414 Lee Xiong, Chenyan Xiong, Ye Li, Kwok-Fung Tang, Jialin Liu, Paul N Bennett, Junaid Ahmed, and Arnold Overwijk. Approximate nearest neighbor negative contrastive learning for dense text retrieval. In *International Conference on Learning Representations*.
- 415 Yingqi Qu, Yuchen Ding, Jing Liu, Kai Liu, Ruiyang Ren, Wayne Xin Zhao, Daxiang Dong, Hua Wu, and Haifeng Wang. Rocketqa: An optimized training approach to dense passage retrieval for open-domain question answering. In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 5835–5847, 2021.
- 416 Liang Wang, Nan Yang, Xiaolong Huang, Binxing Jiao, Linjun Yang, Daxin Jiang, Rangan Majumder, and Furu Wei. Text embeddings by weakly-supervised contrastive pre-training. *arXiv preprint arXiv:2212.03533*, 2022.
- 417 Patrick Lewis, Barlas Oguz, Wenhan Xiong, Fabio Petroni, Scott Yih, and Sebastian Riedel. Boosted dense retriever. In *Proceedings of the 2022 Association for Computational Linguistics: Human Language Technologies*, pages 3102–3117, 2022.
- 418 Liang Wang, Nan Yang, Xiaolong Huang, Binxing Jiao, Linjun Yang, Daxin Jiang, Rangan Majumder, and Furu Wei. Simlm: Pre-training with representation bottleneck for dense passage retrieval. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 2244–2258, 2023.
- 419 Wenhan Xiong, Xiang Lorraine Li, Srini Iyer, Jingfei Du, Patrick S. H. Lewis, William Yang Wang, Yashar Mehdad, Scott Yih, Sebastian Riedel, Douwe Kiela, and Barlas Oguz. Answering complex open-domain questions with multi-hop dense retrieval. In *9th International Conference on Learning Representations, ICLR 2021, Virtual Event, Austria, May 3-7, 2021*, 2021.
- 420 Ori Yoran, Tomer Wolfson, Ori Ram, and Jonathan Berant. Making retrieval-augmented language models robust to irrelevant context. In *The Twelfth International Conference on Learning Representations, ICLR 2024, Vienna, Austria, May 7-11, 2024*. OpenReview.net, 2024.
- 421 Adam Trischler, Tong Wang, Xingdi Yuan, Justin Harris, Alessandro Sordoni, Philip Bachman, and Kaheer Suleman. NewsQA: A machine comprehension dataset. In Phil Blunsom, Antoine Bordes, Kyunghyun Cho, Shay Cohen, Chris Dyer, Edward Grefenstette, Karl Moritz Hermann, Laura Rimell, Jason Weston, and Scott Yih, editors, *Proceedings of the 2nd Workshop on Representation Learning for NLP*, pages 191–200, Vancouver, Canada, August 2017. Association for Computational Linguistics.
- 422 Peter Clark, Isaac Cowhey, Oren Etzioni, Tushar Khot, Ashish Sabharwal, Carissa Schoenick, and Oyvind Tafjord. Think you have solved question answering? try arc, the ai2 reasoning challenge, 2018.
- 423 Pranav Rajpurkar, Robin Jia, and Percy Liang. Know what you don’t know: Unanswerable questions for SQuAD. In Iryna Gurevych and Yusuke Miyao, editors, *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, pages 784–789, Melbourne, Australia, July 2018. Association for Computational Linguistics.
- 424 Yi Yang, Scott Wen-tau Yih, and Chris Meek. Wikiqa: A challenge dataset for open-domain question answering. In *Proceedings of the 2015 Conference on Empirical Methods in Natural Language Processing*. ACL - Association for Computational Linguistics, September 2015.
- 425 Matthew Dunn, Levent Sagun, Mike Higgins, V. Ugur Guney, Volkan Cirik, and Kyunghyun Cho. Searchqa: A new q&a dataset augmented with context from a search engine, 2017.
- 426 Zhilin Yang, Peng Qi, Saizheng Zhang, Joshua Bengio, William Cohen, Ruslan Salakhutdinov, and Christopher D. Manning. HotpotQA: A dataset for diverse, explainable multi-hop question answering. In Ellen Riloff, David Chiang, Julia Hockenmaier, and Jun’ichi Tsujii, editors, *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 2369–2380, Brussels, Belgium, October-November 2018. Association for Computational Linguistics.
- 427 Tom Kwiatkowski, Jennimaria Palomaki, Olivia Redfield, Michael Collins, Ankur Parikh, Chris Alberti, Danielle Epstein, Illia Polosukhin, Jacob Devlin, Kenton Lee, Kristina Toutanova, Llion Jones, Matthew Kelcey, Ming-Wei Chang, Andrew M. Dai, Jakob Uszkoreit, Quoc Le, and Slav Petrov. Natural questions: A benchmark for question answering research. *Transactions of the Association for Computational Linguistics*, 7:452–466, 2019.
- 428 Xanh Ho, Anh-Khoa Duong Nguyen, Saku Sugawara, and Akiko Aizawa. Constructing a multi-hop QA dataset for comprehensive evaluation of reasoning steps. In Donia Scott, Nuria Bel, and Chengqing Zong, editors, *Proceedings of the 28th International Conference on Computational Linguistics*, pages 6609–6625, Barcelona, Spain (Online), December 2020. International Committee on Computational Linguistics.
- 429 James Ferguson, Matt Gardner, Hannaneh Hajishirzi, Tushar Khot, and Pradeep Dasigi. IIRC: A dataset of incomplete information reading comprehension questions. In Bonnie Webber, Trevor Cohn, Yulan He, and Yang Liu, editors, *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 1137–1147, Online, November 2020. Association for Computational Linguistics.
- 430 Andrew Zhu, Alyssa Hwang, Liam Dugan, and Chris Callison-Burch. FanOutQA: A multi-hop, multi-document question answering benchmark for large language models. In Lun-Wei Ku, Andre Martins, and Vivek Srikumar, editors, *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, pages 18–37, Bangkok, Thailand, August 2024. Association for Computational Linguistics.
- 431 Alexander Miller, Adam Fisch, Jesse Dodge, Amir-Hossein Karimi, Antoine Bordes, and Jason Weston. Key-value memory networks for directly reading documents. In *Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing*, pages 1400–1409, 2016.
- 432 Wenhu Chen, Hanwen Zha, Zhiyu Chen, Wenhan Xiong, Hong Wang, and William Yang Wang. Hybridqa: A dataset of multi-hop question answering over tabular and textual data. In *Findings of the Association for Computational Linguistics: EMNLP 2020*, pages 1026–1036, 2020.
- 433 Soumajit Pramanik, Jesujoba Alabi, Rishiraj Saha Roy, and Gerhard Weikum. Uniqrn: unified question answering over rdf knowledge graphs and natural language text. *Journal of Web Semantics*, 83:100833, 2024.
- 434 Shicheng Liu, Jialiang Xu, Wesley Tjangnaka, Sina Semnani, Chen Yu, and Monica Lam. Suql: Conversational search

- over structured and unstructured data with large language models. In *Findings of the Association for Computational Linguistics: NAACL 2024*, pages 4535–4555, 2024.
- 435 Wenhui Chen, Ming-Wei Chang, Eva Schlinger, William Yang Wang, and William W. Cohen. Open question answering over tables and text. In *International Conference on Learning Representations*, 2021.
- 436 Wenting Zhao, Ye Liu, Tong Niu, Yao Wan, Philip Yu, Shafiq Joty, Yingbo Zhou, and Semih Yavuz. DIVKNOWQA: Assessing the reasoning ability of LLMs via open-domain question answering over knowledge base and text. In *Findings of the Association for Computational Linguistics: NAACL 2024*, pages 51–68, 2024.
- 437 Kelvin Han and Claire Gardent. Generating and answering simple and complex questions from text and from knowledge graphs. In *The 13th International Joint Conference on Natural Language Processing and the 3rd Conference of the Asia-Pacific Chapter of the Association for Computational Linguistics (IJCNLP-AAACL 2023)*, 2023.
- 438 Kaixin Ma, Hao Cheng, Xiaodong Liu, Eric Nyberg, and Jianfeng Gao. Open-domain question answering via chain of reasoning over heterogeneous knowledge. In *Findings of the Association for Computational Linguistics: EMNLP 2022*, pages 5360–5374, 2022.
- 439 Philipp Christmann, Rishiraj Saha Roy, and Gerhard Weikum. Compmix: A benchmark for heterogeneous question answering. In *Companion Proceedings of the ACM Web Conference 2024*, WWW '24, page 1091–1094, New York, NY, USA, 2024. Association for Computing Machinery.
- 440 Alon Talmor, Ori Yoran, Amnon Catav, Dan Lahav, Yizhong Wang, Akari Asai, Gabriel Ilharco, Hannaneh Hajishirzi, and Jonathan Berant. Multimodal{qa}: complex question answering over text, tables and images. In *International Conference on Learning Representations*, 2021.
- 441 Darryl Hannan, Akshay Jain, and Mohit Bansal. Manymodalqa: Modality disambiguation and qa over diverse inputs. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 34, pages 7879–7886, 2020.
- 442 Fengbin Zhu, Wenqiang Lei, Youcheng Huang, Chao Wang, Shuo Zhang, Jiancheng Lv, Fuli Feng, and Tat-Seng Chua. TAT-QA: A question answering benchmark on a hybrid of tabular and textual content in finance. In Chengqing Zong, Fei Xia, Wenjie Li, and Roberto Navigli, editors, *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pages 3277–3287, Online, August 2021. Association for Computational Linguistics.
- 443 Zhiyu Chen, Wenhui Chen, Charese Smiley, Sameena Shah, Iana Borova, Dylan Langdon, Reema Moussa, Matt Beane, Ting-Hao Huang, Bryan Routledge, and William Yang Wang. FinQA: A dataset of numerical reasoning over financial data. In Marie-Francine Moens, Xuanjing Huang, Lucia Specia, and Scott Wen-tau Yih, editors, *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 3697–3711, Online and Punta Cana, Dominican Republic, November 2021. Association for Computational Linguistics.
- 444 Xiaoyu Shen, Gianni Barlacchi, Marco Del Tredici, Weiwei Cheng, Bill Byrne, and Adrià Gispert. Product answer generation from heterogeneous sources: A new benchmark and best practices. In Shervin Malmasi, Oleg Rokhlenko, Nicola Ueffing, Ido Guy, Eugene Agichtein, and Surya Kallumadi, editors, *Proceedings of the Fifth Workshop on e-Commerce and NLP (ECNLP 5)*, pages 99–110, Dublin, Ireland, May 2022. Association for Computational Linguistics.
- 445 Zixuan Ke, Yijia Shao, Haowei Lin, Tatsuya Konishi, Gyuhak Kim, and Bing Liu. Continual pre-training of language models. In *The Eleventh International Conference on Learning Representations*, 2023.
- 446 Yichi Zhang, Zhuo Chen, Yin Fang, Yanxi Lu, Li Fangming, Wen Zhang, and HuaJun Chen. Knowledgeable preference alignment for LLMs in domain-specific question answering. In Lun-Wei Ku, Andre Martins, and Vivek Srikumar, editors, *Findings of the Association for Computational Linguistics: ACL 2024*, pages 891–904, Bangkok, Thailand, August 2024. Association for Computational Linguistics.
- 447 Abulhair Saparov and He He. Language models are greedy reasoners: A systematic formal analysis of chain-of-thought, 2023.
- 448 Nouha Dziri, Ximing Lu, Melanie Sclar, Xiang (Lorraine) Li, Liwei Jiang, Bill Yuchen Lin, Sean Welleck, Peter West, Chandra Bhagavatula, Ronan Le Bras, Jena Hwang, Soumya Sanyal, Xiang Ren, Allyson Ettinger, Zaid Harchaoui, and Yejin Choi. Faith and fate: Limits of transformers on compositionality. In A. Oh, T. Naumann, A. Globerson, K. Saenko, M. Hardt, and S. Levine, editors, *Advances in Neural Information Processing Systems*, volume 36, pages 70293–70332. Curran Associates, Inc., 2023.
- 449 Andreas Opedal, Haruki Shirakami, Bernhard Schölkopf, Abulhair Saparov, and Mrinmaya Sachan. Mathgap: Out-of-distribution evaluation on problems with arbitrarily complex proofs, 2024.
- 450 Rachneet Sachdeva, Martin Tutek, and Iryna Gurevych. Catfood: Counterfactual augmented training for improving out-of-domain performance and calibration. In *Proceedings of the 18th Conference of the European Chapter of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1876–1898, 2024.
- 451 Adyasha Maharana and Mohit Bansal. Grada: Graph generative data augmentation for commonsense reasoning. In *Proceedings of the 29th International Conference on Computational Linguistics*, pages 4499–4516, 2022.
- 452 Anji Liu, Hongming Xu, Guy Van den Broeck, and Yitao Liang. Out-of-distribution generalization by neural-symbolic joint training. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 37, pages 12252–12259, 2023.
- 453 Artur SD'Avila Garcez, Luis C Lamb, and Dov M Gabbay. *Neural-symbolic cognitive reasoning*. Springer Science & Business Media, 2008.
- 454 Mayi Xu, Yongqi Li, Ke Sun, and Tiejun Qian. Adaption-of-thought: Learning question difficulty improves large language models for reasoning. In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, pages 5468–5495, 2024.
- 455 Xingyu Chen, Jiahao Xu, Tian Liang, Zhiwei He, Jianhui Pang, Dian Yu, Linfeng Song, Qiuzhi Liu, Mengfei Zhou, Zhuosheng Zhang, Rui Wang, Zhaopeng Tu, Haitao Mi, and Dong Yu. Do not think that much for $2+3=?$ on the overthinking of o1-like llms. *arXiv preprint arXiv:2412.21187*, 2024.
- 456 Fabio Del Vigna12, Andrea Cimino23, Felice Dell'Orletta, Marinella Petrocchi, and Maurizio Tesconi. Hate me, hate me not: Hate speech detection on facebook. In *The first ITASEC*, pages 86–95, 2017.
- 457 Paula Fortuna and Sérgio Nunes. A survey on automatic detection of hate speech in text. *ACM Computing Surveys (CSUR)*, 51(4):1–30, 2018.
- 458 Paul Pu Liang, Chiyu Wu, Louis-Philippe Morency, and Ruslan Salakhutdinov. Towards understanding and mitigating social biases in language models. In *ICML*, pages 6565–6576. PMLR, 2021.
- 459 Janosch Haber, Bertie Vidgen, Matthew Chapman, Vibhor Agarwal, Roy Ka-Wei Lee, Yong Keong Yap, and Paul Röttger. Improving the detection of multilingual online attacks with rich social media data from singapore. In *61st ACL*, pages 12705–12721, 2023.
- 460 Emad Kebriaei, Ali Homayouni, Roghayeh Faraji, Armita Razavi, Azadeh Shakery, Hesham Faili, and Yadollah Yaghoobzadeh. Persian offensive language detection. *Machine Learning*, 113(7):4359–4379, 2024.
- 461 Yunfan Gao, Yun Xiong, Xinyu Gao, Kangxiang Jia, Jinliu Pan, Yuxi Bi, Yi Dai, Jiawei Sun, and Haofen Wang. Retrieval-augmented generation for large language models: A survey. *arXiv preprint arXiv:2312.10997*, 2023.